

CMA-ES and Advanced Adaptation Mechanisms

Youhei Akimoto¹ & Anne Auger² & Nikolaus Hansen²

1. Shinshu University, Nagano, Japan

2. Inria, Research Centre Saclay, France

y.akimoto@shinshu-u.ac.jp
anne.auger@inria.fr
nikolaus.hansen@inria.fr

We are happy to answer questions at any time.

1

2

Problem Statement Black Box Optimization and Its Difficulties

Problem Statement

Continuous Domain Search/Optimization

- Task: minimize an objective function (*fitness* function, *loss* function) in continuous domain

$$f : \mathcal{X} \subseteq \mathbb{R}^n \rightarrow \mathbb{R}, \quad \mathbf{x} \mapsto f(\mathbf{x})$$

- Black Box scenario (direct search scenario)



- ▶ gradients are not available or not useful
- ▶ problem domain specific knowledge is used only within the black box, e.g. within an appropriate encoding

- Search costs: number of function evaluations

① Problem Statement

② Evolution Strategy (ES)

③ Step-Size Adaptation

- ▶ Why Step-Size Control
- ▶ Path Length Control (CSA)
- ▶ Limitations of CSA and its Alternatives

④ Covariance Matrix Adaptation (CMA)

- ▶ Rank-One Update and Cumulation
- ▶ Rank- μ Update
- ▶ Active Covariance Update

⑤ Design Principle

- ▶ Theoretical Foundations
- ▶ Variants for Large Scale Problems

⑥ Summary and Final Remarks

3

533

4

Problem Statement

Continuous Domain Search/Optimization

- Goal

- ▶ fast convergence to the global optimum
- ▶ solution x with small function value $f(x)$ with least search cost
... or to a robust solution x
there are two conflicting objectives

- Typical Examples

- ▶ shape optimization (e.g. using CFD)
 - ▶ model calibration
 - ▶ parameter calibration
- curve fitting, airfoils
biological, physical
controller, plants, images

- Problems

- ▶ exhaustive search is infeasible
- ▶ naive random search takes too long
- ▶ deterministic search is not successful / takes too long

Approach: stochastic search, Evolutionary Algorithms

5

Problem Statement

Continuous Domain Search/Optimization

- Goal

- ▶ fast convergence to the global optimum
- ▶ solution x with small function value $f(x)$ with least search cost
... or to a robust solution x
there are two conflicting objectives

- Typical Examples

- ▶ shape optimization (e.g. using CFD)
 - ▶ model calibration
 - ▶ parameter calibration
- curve fitting, airfoils
biological, physical
controller, plants, images

- Problems

- ▶ exhaustive search is infeasible
- ▶ naive random search takes too long
- ▶ deterministic search is not successful / takes too long

Approach: stochastic search, Evolutionary Algorithms

7

Problem Statement

Continuous Domain Search/Optimization

- Goal

- ▶ fast convergence to the global optimum
- ▶ solution x with small function value $f(x)$ with least search cost
... or to a robust solution x
there are two conflicting objectives

- Typical Examples

- ▶ shape optimization (e.g. using CFD)
 - ▶ model calibration
 - ▶ parameter calibration
- curve fitting, airfoils
biological, physical
controller, plants, images

- Problems

- ▶ exhaustive search is infeasible
- ▶ naive random search takes too long
- ▶ deterministic search is not successful / takes too long

Approach: stochastic search, Evolutionary Algorithms

6

Objective Function Properties

We assume $f : \mathcal{X} \subset \mathbb{R}^n \rightarrow \mathbb{R}$ to be non-linear, non-separable and to have at least moderate dimensionality, say $n \leq 10$.

Additionally, f can be

- non-convex
 - multimodal
 - non-smooth
 - discontinuous, plateaus
 - ill-conditioned
 - noisy
 - ...
- there are possibly many local optima
- derivatives do not exist

Goal : cope with any of these function properties

they are related to real-world problems

Objective Function Properties

We assume $f : \mathcal{X} \subset \mathbb{R}^n \rightarrow \mathbb{R}$ to be *non-linear, non-separable* and to have at least moderate dimensionality, say $n \leq 10$.

Additionally, f can be

- non-convex
 - multimodal
 - non-smooth
 - discontinuous, plateaus
 - ill-conditioned
 - noisy
 - ...
- there are possibly many local optima
derivatives do not exist

Goal: cope with any of these function properties
they are related to real-world problems

9



What Makes a Function Difficult to Solve?

...and what can be done

The Problem	Possible Approaches
Dimensionality	exploiting the problem structure <i>separability, locality/neighborhood, encoding</i>
Ill-conditioning	second order approach <i>changes the neighborhood metric</i>
Ruggedness	<i>non-local</i> policy, large sampling width (step-size) <i>as large as possible while preserving a reasonable convergence speed</i> <i>population-based</i> method, stochastic, non-elitistic recombination operator <i>serves as repair mechanism</i> restarts

11

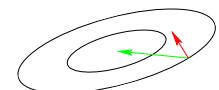
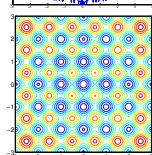
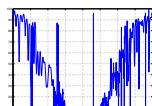
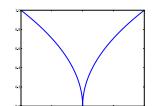
...metaphors



What Makes a Function Difficult to Solve?

Why stochastic search?

- non-linear, non-quadratic, non-convex
on linear and quadratic functions much better search policies are available
- ruggedness
non-smooth, discontinuous, multimodal, and/or noisy function
- dimensionality (size of search space)
(considerably) larger than three
- non-separability
dependencies between the objective variables
- ill-conditioning



10



What Makes a Function Difficult to Solve?

...and what can be done

The Problem	Possible Approaches
Dimensionality	exploiting the problem structure <i>separability, locality/neighborhood, encoding</i>
Ill-conditioning	second order approach <i>changes the neighborhood metric</i>
Ruggedness	<i>non-local</i> policy, large sampling width (step-size) <i>as large as possible while preserving a reasonable convergence speed</i> <i>population-based</i> method, stochastic, non-elitistic recombination operator <i>serves as repair mechanism</i> restarts

...metaphors



What Makes a Function Difficult to Solve?

...and what can be done

The Problem	Possible Approaches
Dimensionality	exploiting the problem structure separability, locality/neighborhood, encoding
III-conditioning	second order approach changes the neighborhood metric
Ruggedness	non-local policy, large sampling width (step-size) as large as possible while preserving a reasonable convergence speed population-based method, stochastic, non-elitistic recombination operator serves as repair mechanism restarts

13

...metaphors

Questions?

14

Stochastic Search

A black box search template to minimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters θ , set population size $\lambda \in \mathbb{N}$
While not terminate

- 1 Sample distribution $P(x|\theta) \rightarrow x_1, \dots, x_\lambda \in \mathbb{R}^n$
- 2 Evaluate x_1, \dots, x_λ on f
- 3 Update parameters $\theta \leftarrow F_\theta(\theta, x_1, \dots, x_\lambda, f(x_1), \dots, f(x_\lambda))$

Everything depends on the definition of P and F_θ

deterministic algorithms are covered as well

In many Evolutionary Algorithms the distribution P is implicitly defined via operators on a population, in particular, selection, recombination and mutation

Natural template for (incremental) *Estimation of Distribution Algorithms*

16

1 Problem Statement

2 Evolution Strategy (ES)

3 Step-Size Adaptation

- ▶ Why Step-Size Control
- ▶ Path Length Control (CSA)
- ▶ Limitations of CSA and its Alternatives

4 Covariance Matrix Adaptation (CMA)

- ▶ Rank-One Update and Cumulation
- ▶ Rank- μ Update
- ▶ Active Covariance Update

5 Design Principle

- ▶ Theoretical Foundations
- ▶ Variants for Large Scale Problems

6 Summary and Final Remarks

15

536

Stochastic Search

A black box search template to minimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters θ , set population size $\lambda \in \mathbb{N}$

While not terminate

- ➊ Sample distribution $P(x|\theta) \rightarrow x_1, \dots, x_\lambda \in \mathbb{R}^n$
- ➋ Evaluate x_1, \dots, x_λ on f
- ➌ Update parameters $\theta \leftarrow F_\theta(\theta, x_1, \dots, x_\lambda, f(x_1), \dots, f(x_\lambda))$

Everything depends on the definition of P and F_θ
deterministic algorithms are covered as well

In many Evolutionary Algorithms the distribution P is implicitly defined via **operators on a population**, in particular, selection, recombination and mutation

Natural template for (incremental) *Estimation of Distribution Algorithms*

17

Stochastic Search

A black box search template to minimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters θ , set population size $\lambda \in \mathbb{N}$

While not terminate

- ➊ Sample distribution $P(x|\theta) \rightarrow x_1, \dots, x_\lambda \in \mathbb{R}^n$
- ➋ Evaluate x_1, \dots, x_λ on f
- ➌ Update parameters $\theta \leftarrow F_\theta(\theta, x_1, \dots, x_\lambda, f(x_1), \dots, f(x_\lambda))$

Everything depends on the definition of P and F_θ
deterministic algorithms are covered as well

In many Evolutionary Algorithms the distribution P is implicitly defined via **operators on a population**, in particular, selection, recombination and mutation

Natural template for (incremental) *Estimation of Distribution Algorithms*

18

Stochastic Search

A black box search template to minimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters θ , set population size $\lambda \in \mathbb{N}$

While not terminate

- ➊ Sample distribution $P(x|\theta) \rightarrow x_1, \dots, x_\lambda \in \mathbb{R}^n$
- ➋ Evaluate x_1, \dots, x_λ on f
- ➌ Update parameters $\theta \leftarrow F_\theta(\theta, x_1, \dots, x_\lambda, f(x_1), \dots, f(x_\lambda))$

Everything depends on the definition of P and F_θ
deterministic algorithms are covered as well

In many Evolutionary Algorithms the distribution P is implicitly defined via **operators on a population**, in particular, selection, recombination and mutation

Natural template for (incremental) *Estimation of Distribution Algorithms*

19

Stochastic Search

A black box search template to minimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters θ , set population size $\lambda \in \mathbb{N}$

While not terminate

- ➊ Sample distribution $P(x|\theta) \rightarrow x_1, \dots, x_\lambda \in \mathbb{R}^n$
- ➋ Evaluate x_1, \dots, x_λ on f
- ➌ Update parameters $\theta \leftarrow F_\theta(\theta, x_1, \dots, x_\lambda, f(x_1), \dots, f(x_\lambda))$

Everything depends on the definition of P and F_θ
deterministic algorithms are covered as well

In many Evolutionary Algorithms the distribution P is implicitly defined via **operators on a population**, in particular, selection, recombination and mutation

Natural template for (incremental) *Estimation of Distribution Algorithms*

20

Stochastic Search

A black box search template to minimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters θ , set population size $\lambda \in \mathbb{N}$

While not terminate

- ➊ Sample distribution $P(x|\theta) \rightarrow x_1, \dots, x_\lambda \in \mathbb{R}^n$
- ➋ Evaluate x_1, \dots, x_λ on f
- ➌ Update parameters $\theta \leftarrow F_\theta(\theta, x_1, \dots, x_\lambda, f(x_1), \dots, f(x_\lambda))$

Everything depends on the definition of P and F_θ
deterministic algorithms are covered as well

In many Evolutionary Algorithms the distribution P is implicitly defined via **operators on a population**, in particular, selection, recombination and mutation

Natural template for (incremental) *Estimation of Distribution Algorithms* ↗
21

Stochastic Search

A black box search template to minimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters θ , set population size $\lambda \in \mathbb{N}$

While not terminate

- ➊ Sample distribution $P(x|\theta) \rightarrow x_1, \dots, x_\lambda \in \mathbb{R}^n$
- ➋ Evaluate x_1, \dots, x_λ on f
- ➌ Update parameters $\theta \leftarrow F_\theta(\theta, x_1, \dots, x_\lambda, f(x_1), \dots, f(x_\lambda))$

Everything depends on the definition of P and F_θ
deterministic algorithms are covered as well

In many Evolutionary Algorithms the distribution P is implicitly defined via **operators on a population**, in particular, selection, recombination and mutation

Natural template for (incremental) *Estimation of Distribution Algorithms* ↗
23

Stochastic Search

A black box search template to minimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters θ , set population size $\lambda \in \mathbb{N}$

While not terminate

- ➊ Sample distribution $P(x|\theta) \rightarrow x_1, \dots, x_\lambda \in \mathbb{R}^n$
- ➋ Evaluate x_1, \dots, x_λ on f
- ➌ Update parameters $\theta \leftarrow F_\theta(\theta, x_1, \dots, x_\lambda, f(x_1), \dots, f(x_\lambda))$

Everything depends on the definition of P and F_θ
deterministic algorithms are covered as well

In many Evolutionary Algorithms the distribution P is implicitly defined via **operators on a population**, in particular, selection, recombination and mutation

Natural template for (incremental) *Estimation of Distribution Algorithms* ↗
22

The CMA-ES

Input: $m \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, λ

Initialize: $C = I$, and $p_c = 0$, $p_\sigma = 0$,

Set: $c_e \approx 4/n$, $c_\sigma \approx 4/n$, $c_1 \approx 2/n^2$, $c_\mu \approx \mu_w/n^2$, $c_1 + c_\mu \leq 1$, $d_\sigma \approx 1 + \sqrt{\frac{\mu_w}{n}}$, and $w_{i=1..n}$ such that $\mu_w = \frac{1}{\sum_{i=1}^n w_i^2} \approx 0.3 \lambda$

While not terminate

$x_i = m + \sigma y_i$, $y_i \sim \mathcal{N}_i(\mathbf{0}, C)$, for $i = 1, \dots, \lambda$

$m \leftarrow \sum_{i=1}^\lambda w_i x_{i:\lambda} = m + \sigma y_w$ where $y_w = \sum_{i=1}^\lambda w_i y_{i:\lambda}$

$p_c \leftarrow (1 - c_e) p_c + \mathbf{1}_{\{\|p_c\| < 1.5\sqrt{n}\}} \sqrt{1 - (1 - c_e)^2} \sqrt{\mu_w} y_w$

$p_\sigma \leftarrow (1 - c_\sigma) p_\sigma + \sqrt{1 - (1 - c_\sigma)^2} \sqrt{\mu_w} C^{-\frac{1}{2}} y_w$

$C \leftarrow (1 - c_1 - c_\mu) C + c_1 p_c p_c^T + c_\mu \sum_{i=1}^\lambda w_i y_{i:\lambda} y_{i:\lambda}^T$

$\sigma \leftarrow \sigma \times \exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|p_\sigma\|}{E\|\mathcal{N}(0, I)\|} - 1 \right) \right)$

Not covered on this slide: termination, restarts, useful output, boundaries and encoding

Evolution Strategies

New search points are sampled normally distributed

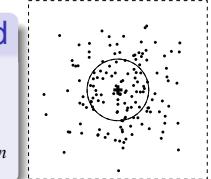
$$\mathbf{x}_i \sim \mathbf{m} + \sigma \mathcal{N}(\mathbf{0}, \mathbf{C}) \quad \text{for } i = 1, \dots, \lambda$$

as perturbations of \mathbf{m} , where $\mathbf{x}_i, \mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, $\mathbf{C} \in \mathbb{R}^{n \times n}$
where

- the mean vector $\mathbf{m} \in \mathbb{R}^n$ represents the favorite solution
- the so-called step-size $\sigma \in \mathbb{R}_+$ controls the step length
- the covariance matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the shape of the distribution ellipsoid

here, all new points are sampled with the same parameters

The question remains how to update \mathbf{m} , \mathbf{C} , and σ .



25

Evolution Strategies

New search points are sampled normally distributed

$$\mathbf{x}_i \sim \mathbf{m} + \sigma \mathcal{N}(\mathbf{0}, \mathbf{C}) \quad \text{for } i = 1, \dots, \lambda$$

as perturbations of \mathbf{m} , where $\mathbf{x}_i, \mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, $\mathbf{C} \in \mathbb{R}^{n \times n}$

where

- the mean vector $\mathbf{m} \in \mathbb{R}^n$ represents the favorite solution
- the so-called step-size $\sigma \in \mathbb{R}_+$ controls the step length
- the covariance matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the shape of the distribution ellipsoid

here, all new points are sampled with the same parameters

The question remains how to update \mathbf{m} , \mathbf{C} , and σ .

26

Evolution Strategies

Terminology

Let μ : # of parents, λ : # of offspring

Plus (elitist) and comma (non-elitist) selection

$(\mu + \lambda)$ -ES: selection in {parents} \cup {offspring}

(μ, λ) -ES: selection in {offspring}

$(1 + 1)$ -ES

Sample one offspring from parent \mathbf{m}

$$\mathbf{x} = \mathbf{m} + \sigma \mathcal{N}(\mathbf{0}, \mathbf{C})$$

If \mathbf{x} better than \mathbf{m} select

$$\mathbf{m} \leftarrow \mathbf{x}$$

The $(\mu/\mu, \lambda)$ -ES

Non-elitist selection and intermediate (weighted) recombination

Given the i -th solution point $\mathbf{x}_i = \mathbf{m} + \sigma \underbrace{\mathcal{N}_i(\mathbf{0}, \mathbf{C})}_{=: \mathbf{y}_i} = \mathbf{m} + \sigma \mathbf{y}_i$

Let $\mathbf{x}_{i:\lambda}$ the i -th ranked solution point, such that $f(\mathbf{x}_{1:\lambda}) \leq \dots \leq f(\mathbf{x}_{\lambda:\lambda})$.

The new mean reads

$$\mathbf{m} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}$$

where

$$w_1 \geq \dots \geq w_{\mu} > 0, \quad \sum_{i=1}^{\mu} w_i = 1, \quad \frac{1}{\sum_{i=1}^{\mu} w_i^2} =: \mu_w \approx \frac{\lambda}{4}$$

The best μ points are selected from the new solutions (non-elitistic) and weighted intermediate recombination is applied.

29

The $(\mu/\mu, \lambda)$ -ES

Non-elitist selection and intermediate (weighted) recombination

Given the i -th solution point $\mathbf{x}_i = \mathbf{m} + \sigma \underbrace{\mathcal{N}_i(\mathbf{0}, \mathbf{C})}_{=: \mathbf{y}_i} = \mathbf{m} + \sigma \mathbf{y}_i$

Let $\mathbf{x}_{i:\lambda}$ the i -th ranked solution point, such that $f(\mathbf{x}_{1:\lambda}) \leq \dots \leq f(\mathbf{x}_{\lambda:\lambda})$.

The new mean reads

$$\mathbf{m} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}$$

where

$$w_1 \geq \dots \geq w_{\mu} > 0, \quad \sum_{i=1}^{\mu} w_i = 1, \quad \frac{1}{\sum_{i=1}^{\mu} w_i^2} =: \mu_w \approx \frac{\lambda}{4}$$

The best μ points are selected from the new solutions (non-elitistic) and weighted intermediate recombination is applied.

30

The $(\mu/\mu, \lambda)$ -ES

Non-elitist selection and intermediate (weighted) recombination

Given the i -th solution point $\mathbf{x}_i = \mathbf{m} + \sigma \underbrace{\mathcal{N}_i(\mathbf{0}, \mathbf{C})}_{=: \mathbf{y}_i} = \mathbf{m} + \sigma \mathbf{y}_i$

Let $\mathbf{x}_{i:\lambda}$ the i -th ranked solution point, such that $f(\mathbf{x}_{1:\lambda}) \leq \dots \leq f(\mathbf{x}_{\lambda:\lambda})$.

The new mean reads

$$\mathbf{m} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda} = \mathbf{m} + \sigma \sum_{i=1}^{\mu} w_i \underbrace{\mathbf{y}_{i:\lambda}}_{=: \mathbf{y}_w}$$

where

$$w_1 \geq \dots \geq w_{\mu} > 0, \quad \sum_{i=1}^{\mu} w_i = 1, \quad \frac{1}{\sum_{i=1}^{\mu} w_i^2} =: \mu_w \approx \frac{\lambda}{4}$$

The best μ points are selected from the new solutions (non-elitistic) and weighted intermediate recombination is applied.

31

Evolution Strategies

Recalling

New search points are sampled normally distributed

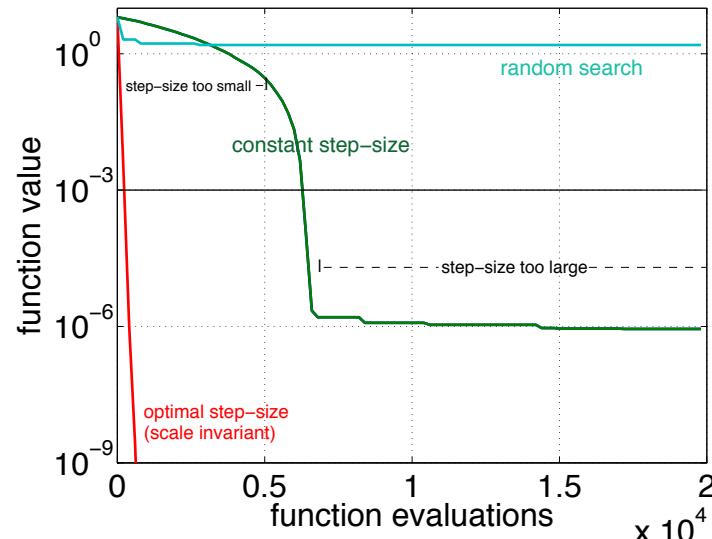
$$\mathbf{x}_i \sim \mathbf{m} + \sigma \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \quad \text{for } i = 1, \dots, \lambda$$

as perturbations of \mathbf{m} , where $\mathbf{x}_i, \mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, $\mathbf{C} \in \mathbb{R}^{n \times n}$
where

- the mean vector $\mathbf{m} \in \mathbb{R}^n$ represents the favorite solution
- the so-called step-size $\sigma \in \mathbb{R}_+$ controls the step length
- the covariance matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the shape of the distribution ellipsoid

The remaining question is how to update \mathbf{C} .

Why Step-Size Control?



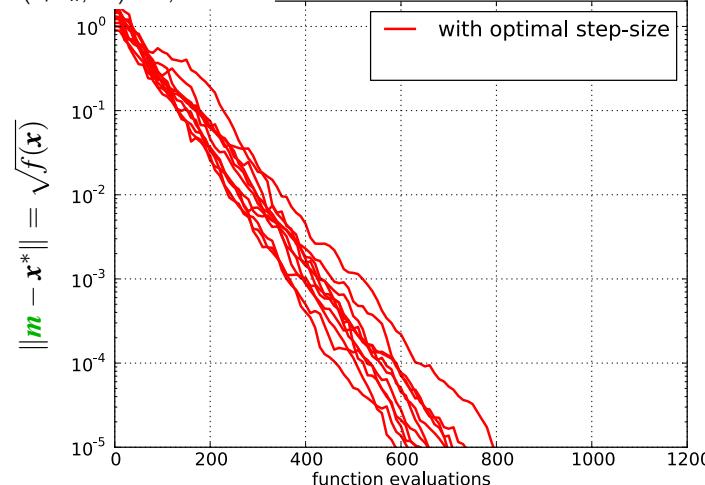
33

(1+1)-ES
(red & green)

$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

in $[-2.2, 0.8]^n$
for $n = 10$

Why Step-Size Control?

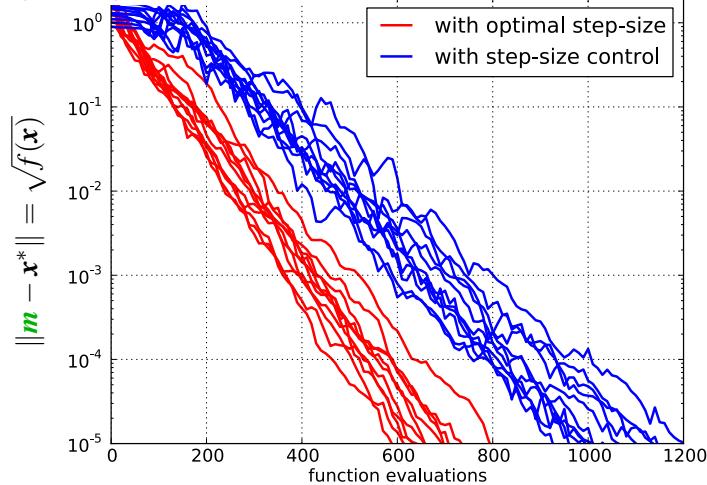
(5/5_w, 10)-ES, 11 runswith optimal step-size σ

$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

for $n = 10$ and
 $\mathbf{x}^0 \in [-0.2, 0.8]^n$

34

Why Step-Size Control?

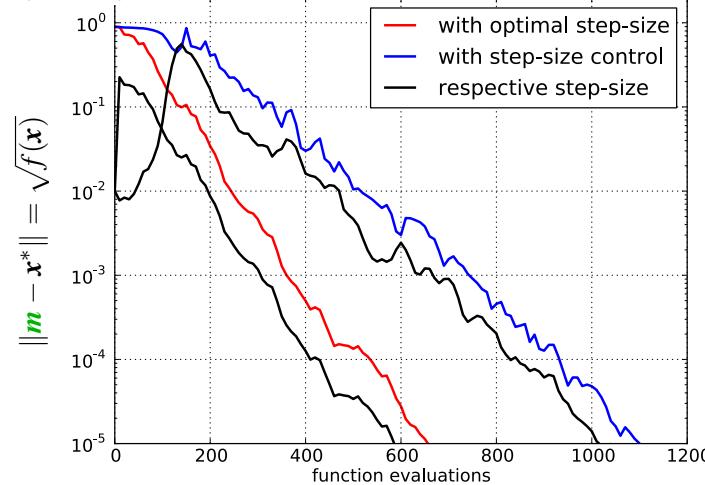
(5/5_w, 10)-ES, 2 times 11 runswith optimal versus adaptive step-size σ with too small initial σ

35

$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

for $n = 10$ and
 $\mathbf{x}^0 \in [-0.2, 0.8]^n$

Why Step-Size Control?

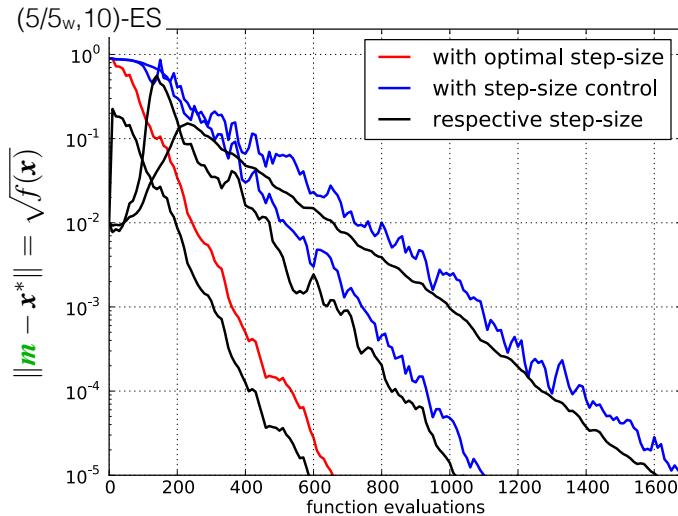
(5/5_w, 10)-EScomparing number of f -evals to reach $\|\mathbf{m}\| = 10^{-5}$: $\frac{1100-100}{650} \approx 1.5$

36

$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

for $n = 10$ and
 $\mathbf{x}^0 \in [-0.2, 0.8]^n$

Why Step-Size Control?



comparing optimal versus default damping parameter $d_\sigma: \frac{1700}{1100} \approx 1.5$

37

Methods for Step-Size Control

- **1/5-th success rule^{a,b}**, often applied with "+"-selection
increase step-size if more than 20% of the new solutions are successful, decrease otherwise
- **σ -self-adaptation^c**, applied with ","-selection
mutation is applied to the step-size and the better, according to the objective function value, is selected
simplified "global" self-adaptation
- **path length control^d** (Cumulative Step-size Adaptation, CSA)^e
self-adaptation derandomized and non-localized

^aRechenberg 1973, *Evolutionsstrategie, Optimierung technischer Systeme nach Prinzipien der biologischen Evolution*, Frommann-Holzboog

^bSchumer and Steiglitz 1968. Adaptive step size random search. *IEEE TAC*

^cSchwefel 1981, *Numerical Optimization of Computer Models*, Wiley

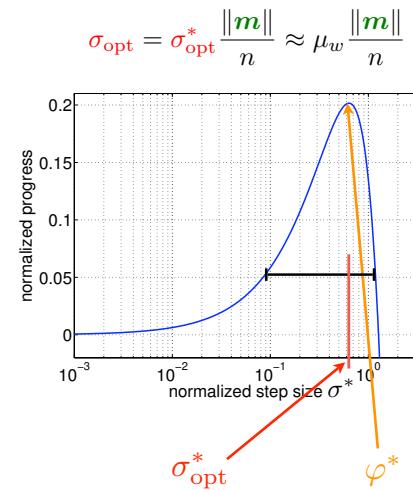
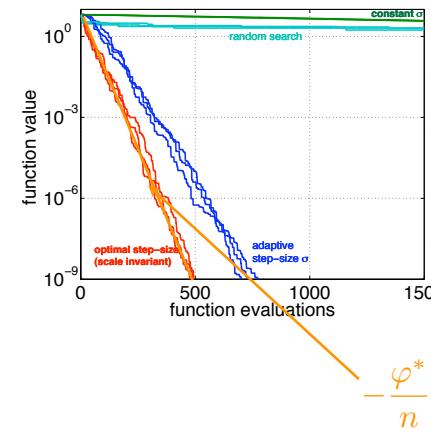
^dHansen & Ostermeier 2001, Completely Derandomized Self-Adaptation in Evolution Strategies, *Evol. Comput.*

9(2)

^eOstermeier et al 1994, Step-size adaptation based on non-local use of selection information, *PPSN IV*

39

Why Step-Size Control?



evolution window refers to the step-size interval (—) where reasonable performance is observed

38

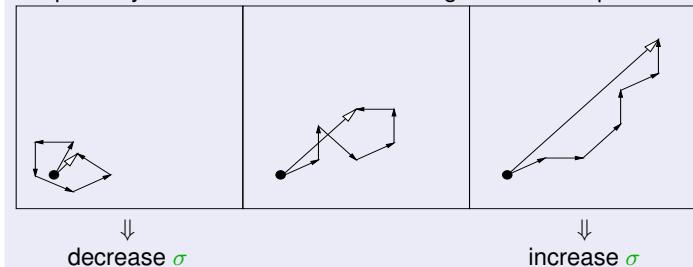
Path Length Control (CSA)

The Concept of Cumulative Step-Size Adaptation

$$\begin{aligned} \mathbf{x}_i &= \mathbf{m} + \sigma \mathbf{y}_i \\ \mathbf{m} &\leftarrow \mathbf{m} + \sigma \mathbf{y}_w \end{aligned}$$

Measure the length of the evolution path

the pathway of the mean vector \mathbf{m} in the generation sequence



loosely speaking steps are

- perpendicular under random selection (in expectation)
- perpendicular in the desired situation (to be most efficient)

40

Path Length Control (CSA)

The Equations

Initialize $\mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, evolution path $\mathbf{p}_\sigma = \mathbf{0}$,
set $c_\sigma \approx 4/n$, $d_\sigma \approx 1$.

$$\begin{aligned}\mathbf{m} &\leftarrow \mathbf{m} + \sigma \mathbf{y}_w \quad \text{where } \mathbf{y}_w = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} && \text{update mean} \\ \mathbf{p}_\sigma &\leftarrow (1 - c_\sigma) \mathbf{p}_\sigma + \underbrace{\sqrt{1 - (1 - c_\sigma)^2}}_{\text{accounts for } 1 - c_\sigma} \underbrace{\sqrt{\mu_w}}_{\text{accounts for } w_i} \mathbf{y}_w \\ \sigma &\leftarrow \sigma \times \exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right) && \text{update step-size} \\ &> 1 \iff \|\mathbf{p}_\sigma\| \text{ is greater than its expectation}\end{aligned}$$

41

Path Length Control (CSA)

The Equations

Initialize $\mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, evolution path $\mathbf{p}_\sigma = \mathbf{0}$,
set $c_\sigma \approx 4/n$, $d_\sigma \approx 1$.

$$\begin{aligned}\mathbf{m} &\leftarrow \mathbf{m} + \sigma \mathbf{y}_w \quad \text{where } \mathbf{y}_w = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} && \text{update mean} \\ \mathbf{p}_\sigma &\leftarrow (1 - c_\sigma) \mathbf{p}_\sigma + \underbrace{\sqrt{1 - (1 - c_\sigma)^2}}_{\text{accounts for } 1 - c_\sigma} \underbrace{\sqrt{\mu_w}}_{\text{accounts for } w_i} \mathbf{y}_w \\ \sigma &\leftarrow \sigma \times \exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right) && \text{update step-size} \\ &> 1 \iff \|\mathbf{p}_\sigma\| \text{ is greater than its expectation}\end{aligned}$$

42

Known Issues of CSA I: With Ineffective Axes

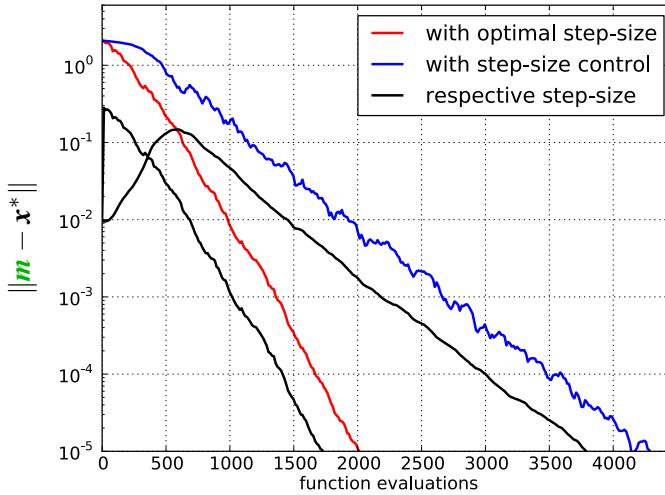
On a function with ineffective axes

- $f(\mathbf{x}) = \sum_{i=1}^M [x]_i^2$, $\mathbf{x} \in \mathbb{R}^N$, $M \leq N$.
- $N - M$ variables do not affect the function value

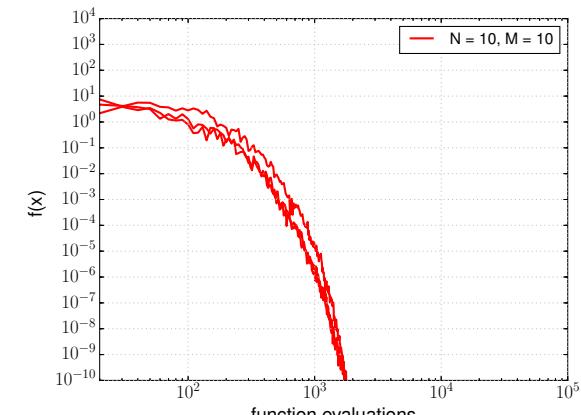
$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

in $[-0.2, 0.8]^n$
for $n = 30$

(5/5, 10)-CSA-ES, default parameters



43

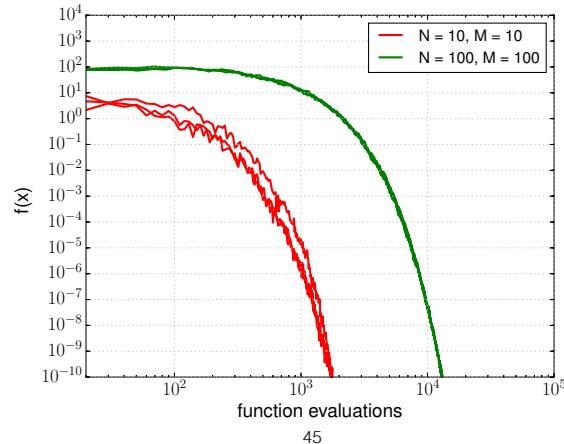


44

Known Issues of CSA I: With Ineffective Axes

On a function with ineffective axes

- $f(\mathbf{x}) = \sum_{i=1}^M [x]_i^2, \quad \mathbf{x} \in \mathbb{R}^N, \quad M \leq N.$
- $N - M$ variables do not affect the function value

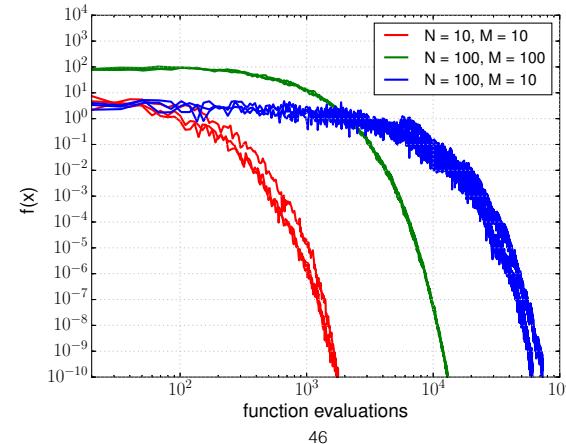


45

Known Issues of CSA I: With Ineffective Axes

On a function with ineffective axes

- $f(\mathbf{x}) = \sum_{i=1}^M [x]_i^2, \quad \mathbf{x} \in \mathbb{R}^N, \quad M \leq N.$
- $N - M$ variables do not affect the function value



46

Known Issues of CSA I: With Ineffective Axes

On a function with ineffective axes

- $f(\mathbf{x}) = \sum_{i=1}^M [x]_i^2, \quad \mathbf{x} \in \mathbb{R}^N, \quad M \leq N.$
- $N - M$ variables do not affect the function value

If $M \ll N$,

- The $N - M$ components of \mathbf{p}_σ are normally distributed
- The signal to change the step size is in the M components

$$\begin{aligned}\mathbf{p}_\sigma &\leftarrow (1 - c_\sigma)\mathbf{p}_\sigma + \sqrt{c_\sigma(2 - c_\sigma)}\sqrt{\mu_w} \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \\ \sigma &\leftarrow \sigma \times \exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)\end{aligned}$$

47

544

Known Issues of CSA II: With a Large Population

With large λ and μ

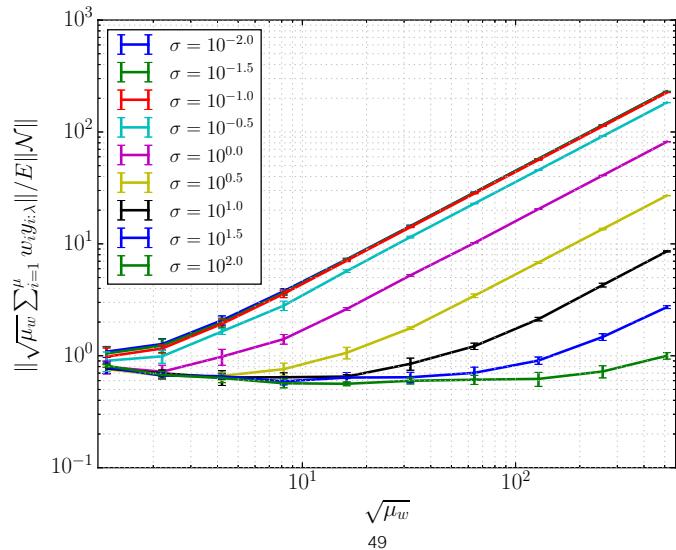
$$\begin{aligned}\mathbf{p}_\sigma &\leftarrow (1 - c_\sigma)\mathbf{p}_\sigma + \sqrt{c_\sigma(2 - c_\sigma)} \underbrace{\sqrt{\mu_w} \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda}}_{\propto \mu \rightarrow \infty \text{ as } \lambda, \mu \rightarrow \infty} \\ \sigma &\leftarrow \sigma \times \exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)\end{aligned}$$

- $\|\mathbf{p}_\sigma\| \propto \sqrt{\mu_w} \implies \frac{\|\mathbf{p}_\sigma\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \propto \sqrt{\frac{\mu_w}{N}} - 1$ unless $\|\mathbf{m} - \mathbf{x}^*\| \ll \sigma$
- $d_\sigma \propto \sqrt{\frac{N}{\mu_w}}$ to prevent a increase of σ in one step, however, the convergence will be very slow.

48

Known Issues of CSA II: With a Large Population

Sphere function ($n = 10$, $\mathbf{m} = (1, 0, \dots, 0)$)



49

Two-Point Step-Size Adaptation (TPA)

- Sample a pair of symmetric point along the previous mean shift

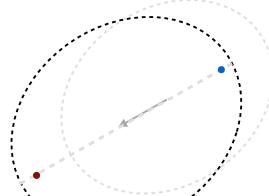
$$\mathbf{x}_{1/2} = \mathbf{m}^{(g)} \pm \sigma^{(g)} \frac{\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|}{\|\mathbf{m}^{(g)} - \mathbf{m}^{(g-1)}\|_{\mathbf{C}^{(g)}}} (\mathbf{m}^{(g)} - \mathbf{m}^{(g-1)}) \quad \|\mathbf{x}\|_{\mathbf{C}} := \mathbf{x}^T \mathbf{C}^{-1} \mathbf{x}$$

- Compare the ranking of x_1 and x_2 among λ current populations

$$s^{(g+1)} = (1 - c_s)s^{(g)} + c_s \underbrace{\frac{\text{rank}(x_2) - \text{rank}(x_1)}{\lambda - 1}}_{>0 \text{ if the previous step still produces a promising solution}}$$

- Update the step-size

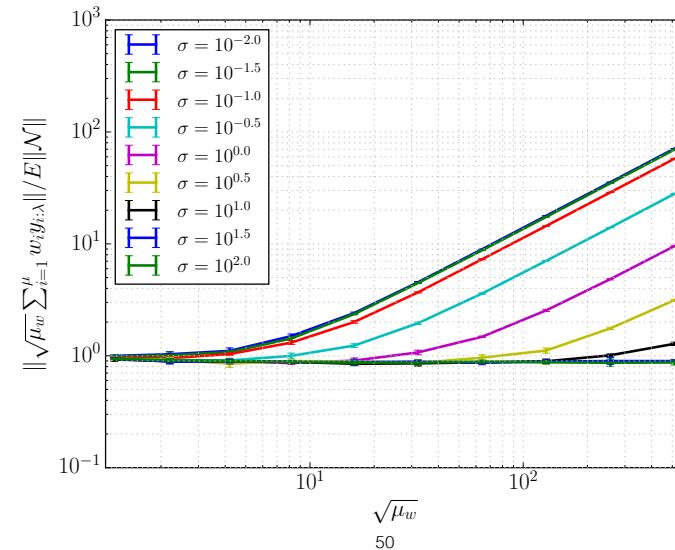
$$\sigma^{(g+1)} = \sigma^{(g)} \exp \left(\frac{s^{(g+1)}}{d_\sigma} \right)$$



[Hansen, 2008] Hansen, N. (2008). CMA-ES with two-point step-size adaptation. [research report] rr-6527, 2008. Inria-00276854v5.
 [Hansen et al., 2014] Hansen, N., Atamna, A., and Auger, A. (2014). How to assess step-size adaptation mechanisms in randomised search. In Parallel Problem Solving from Nature—PPSN XIII, pages 60–69. Springer.

Known Issues of CSA II: With a Large Population

Sphere function ($n = 100$, $\mathbf{m} = (1, 0, \dots, 0)$)



50

Two-Point Step-Size Adaptation (TPA)

- Sample a pair of symmetric point along the previous mean shift

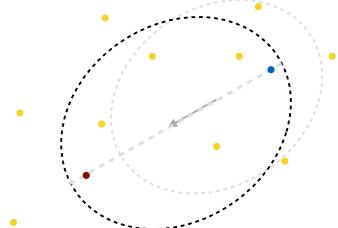
$$\mathbf{x}_{1/2} = \mathbf{m}^{(g)} \pm \sigma^{(g)} \frac{\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|}{\|\mathbf{m}^{(g)} - \mathbf{m}^{(g-1)}\|_{\mathbf{C}^{(g)}}} (\mathbf{m}^{(g)} - \mathbf{m}^{(g-1)}) \quad \|\mathbf{x}\|_{\mathbf{C}} := \mathbf{x}^T \mathbf{C}^{-1} \mathbf{x}$$

- Compare the ranking of x_1 and x_2 among λ current populations

$$s^{(g+1)} = (1 - c_s)s^{(g)} + c_s \underbrace{\frac{\text{rank}(x_2) - \text{rank}(x_1)}{\lambda - 1}}_{>0 \text{ if the previous step still produces a promising solution}}$$

- Update the step-size

$$\sigma^{(g+1)} = \sigma^{(g)} \exp \left(\frac{s^{(g+1)}}{d_\sigma} \right)$$



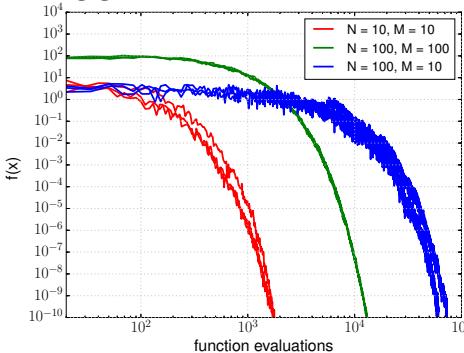
[Hansen, 2008] Hansen, N. (2008). CMA-ES with two-point step-size adaptation. [research report] rr-6527, 2008. Inria-00276854v5.
 [Hansen et al., 2014] Hansen, N., Atamna, A., and Auger, A. (2014). How to assess step-size adaptation mechanisms in randomised search. In Parallel Problem Solving from Nature—PPSN XIII, pages 60–69. Springer.

On Sphere with Ineffective Axes

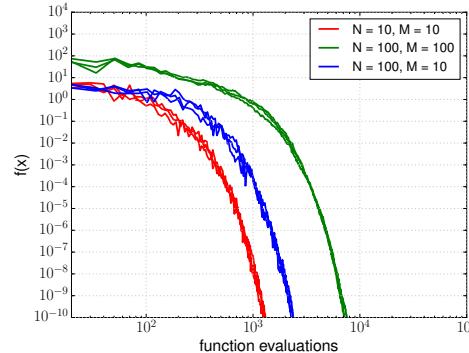
On a function with ineffective axes

- $f(\mathbf{x}) = \sum_{i=1}^M [x]_i^2, \quad \mathbf{x} \in \mathbb{R}^N, \quad M \leq N.$
- $N - M$ variables do not affect the function value

CSA



TPA



53

Alternatives: Success-Based Step-Size Control

comparing the fitness distributions of current and previous iterations

Generalizations of 1/5th-success-rule for non-elitist and multi-recombinant ES

- Median Success Rule [Ait Elhara et al., 2013]
- Population Success Rule [Loshchilov, 2014]

controls a *success probability*

An advantage over CSA and TPA: Cheap Computation

- It depends only on λ .
- cf. CSA and TPA require a computation of $\mathbf{C}^{-1/2}\mathbf{x}$ and $\mathbf{C}^{-1}\mathbf{x}$, respectively.

[Ait Elhara et al., 2013] Ait Elhara, O., Auger, A., and Hansen, N. (2013). A median success rule for non- elitist evolution strategies: Study of feasibility. In Proc. of the GECCO, pages 415–422.

[Loshchilov, 2014] Loshchilov, I. (2014). A computationally efficient limited memory cma-es for large scale optimization. In Proc. of the GECCO, pages 397–404.

Alternatives: Success-Based Step-Size Control

comparing the fitness distributions of current and previous iterations

Generalizations of 1/5th-success-rule for non-elitist and multi-recombinant ES

- Median Success Rule [Ait Elhara et al., 2013]

- Population Success Rule [Loshchilov, 2014]

controls a *success probability*

An advantage over CSA and TPA: Cheap Computation

- It depends only on λ .
- cf. CSA and TPA require a computation of $\mathbf{C}^{-1/2}\mathbf{x}$ and $\mathbf{C}^{-1}\mathbf{x}$, respectively.

[Ait Elhara et al., 2013] Ait Elhara, O., Auger, A., and Hansen, N. (2013). A median success rule for non- elitist evolution strategies: Study of feasibility. In Proc. of the GECCO, pages 415–422.

[Loshchilov, 2014] Loshchilov, I. (2014). A computationally efficient limited memory cma-es for large scale optimization. In Proc. of the GECCO, pages 397–404.

54

- ① Problem Statement
- ② Evolution Strategy (ES)
- ③ Step-Size Adaptation
 - ▶ Why Step-Size Control
 - ▶ Path Length Control (CSA)
 - ▶ Limitations of CSA and its Alternatives
- ④ Covariance Matrix Adaptation (CMA)
 - ▶ Rank-One Update and Cumulation
 - ▶ Rank- μ Update
 - ▶ Active Covariance Update
- ⑤ Design Principle
 - ▶ Theoretical Foundations
 - ▶ Variants for Large Scale Problems
- ⑥ Summary and Final Remarks

Evolution Strategies

Recalling

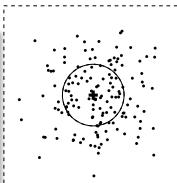
New search points are sampled normally distributed

$$\mathbf{x}_i \sim \mathbf{m} + \sigma \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \quad \text{for } i = 1, \dots, \lambda$$

as perturbations of \mathbf{m} , where $\mathbf{x}_i, \mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, $\mathbf{C} \in \mathbb{R}^{n \times n}$

where

- the mean vector $\mathbf{m} \in \mathbb{R}^n$ represents the favorite solution and $\mathbf{m} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}$
- the so-called step-size $\sigma \in \mathbb{R}_+$ controls the step length
- the covariance matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the shape of the distribution ellipsoid



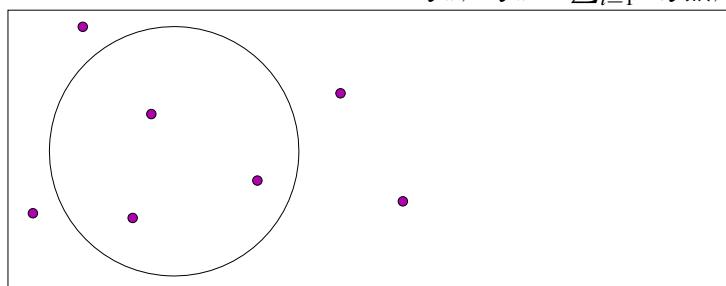
The remaining question is how to update σ and \mathbf{C} .

57

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w, \quad \mathbf{y}_w = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda}, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



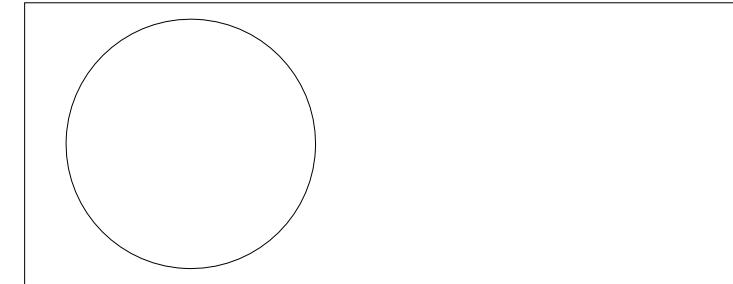
initial distribution, $\mathbf{C} = \mathbf{I}$

59

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w, \quad \mathbf{y}_w = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda}, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



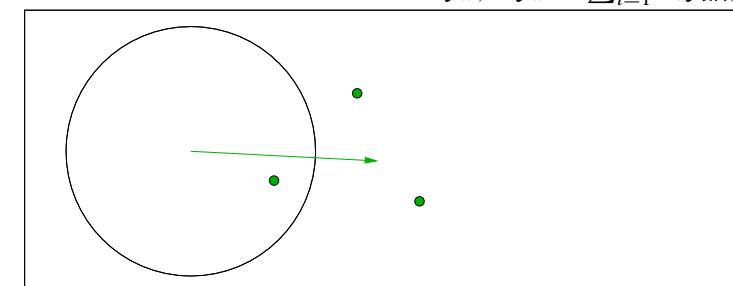
initial distribution, $\mathbf{C} = \mathbf{I}$

58

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w, \quad \mathbf{y}_w = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda}, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



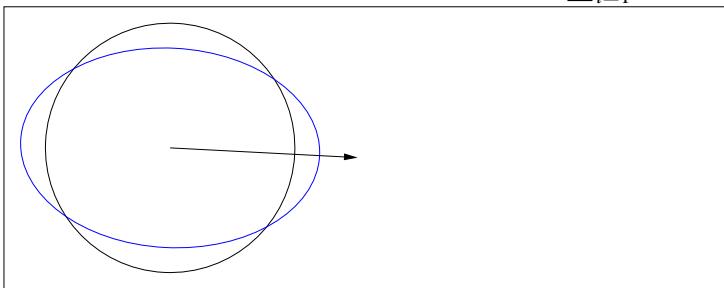
\mathbf{y}_w , movement of the population mean \mathbf{m} (disregarding σ)

60

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w, \quad \mathbf{y}_w = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda}, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



mixture of distribution \mathbf{C} and step \mathbf{y}_w ,

$$\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^T$$

61

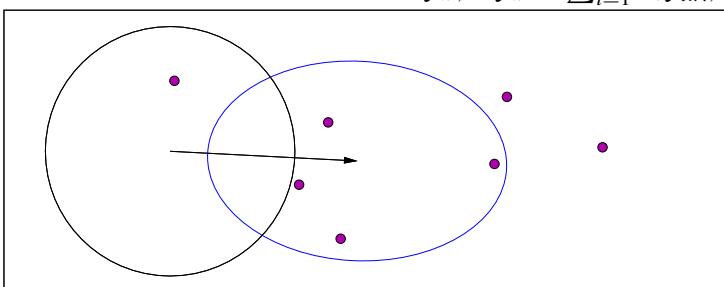
...equations



Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w, \quad \mathbf{y}_w = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda}, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



new distribution (disregarding σ)

63

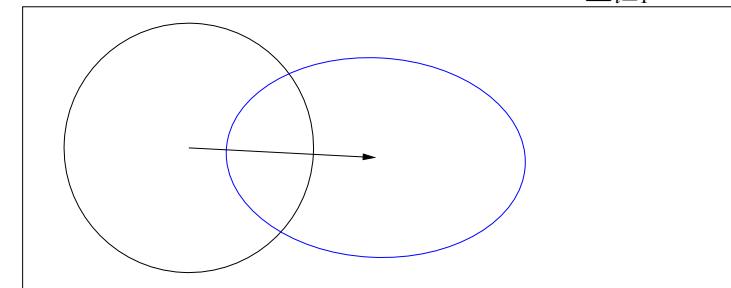
...equations



Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w, \quad \mathbf{y}_w = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda}, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



new distribution (disregarding σ)

62

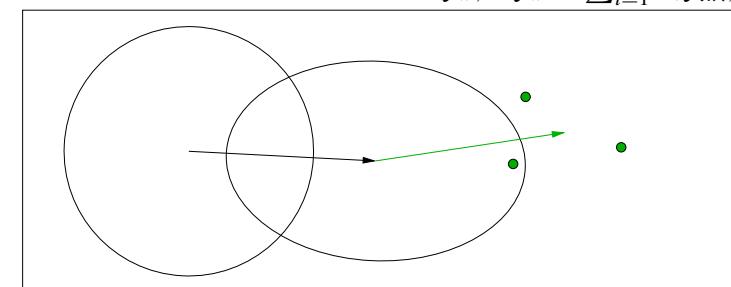
...equations



Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w, \quad \mathbf{y}_w = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda}, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



movement of the population mean \mathbf{m}

64

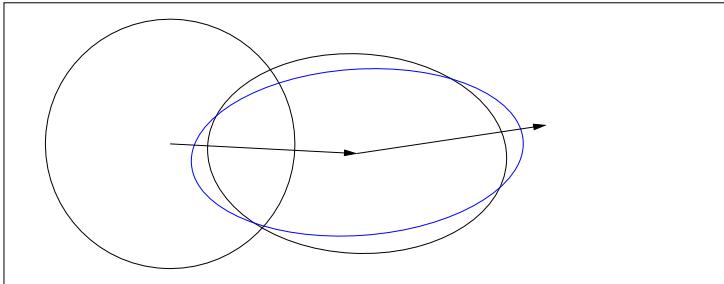
...equations



Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w, \quad \mathbf{y}_w = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda}, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



mixture of distribution \mathbf{C} and step \mathbf{y}_w ,

$$\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^T$$

... equations

65



Covariance Matrix Adaptation

Rank-One Update

Initialize $\mathbf{m} \in \mathbb{R}^n$, and $\mathbf{C} = \mathbf{I}$, set $\sigma = 1$, learning rate $c_{\text{cov}} \approx 2/n^2$

While not terminate

$$\mathbf{x}_i = \mathbf{m} + \sigma \mathbf{y}_i, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C}),$$

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w \quad \text{where } \mathbf{y}_w = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda}$$

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mu_w \underbrace{\mathbf{y}_w \mathbf{y}_w^T}_{\text{rank-one}} \quad \text{where } \mu_w = \frac{1}{\sum_{i=1}^{\mu} \mathbf{w}_i^2} \geq 1$$

The rank-one update has been found independently in several domains^{6 7 8 9}

⁶ Kjellström & Taxén 1981. Stochastic Optimization in System Design, IEEE TCS

⁷ Hansen & Ostermeier 1996. Adapting arbitrary normal mutation distributions in evolution strategies: The covariance matrix adaptation, ICEC

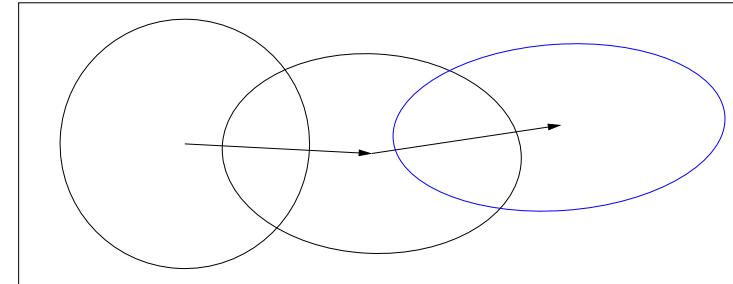
⁸ Ljung 1999. System Identification: Theory for the User

⁹ Haario et al 2001. An adaptive Metropolis algorithm, JSTOR

Covariance Matrix Adaptation

Rank-One Update

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w, \quad \mathbf{y}_w = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda}, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$



new distribution,

$$\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^T$$

the ruling principle: the adaptation increases the likelihood of successful steps, \mathbf{y}_w , to appear again
another viewpoint: the adaptation follows a natural gradient approximation of the expected fitness

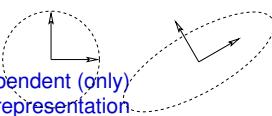


66

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mu_w \mathbf{y}_w \mathbf{y}_w^T$$

covariance matrix adaptation

- learns all pairwise dependencies between variables
off-diagonal entries in the covariance matrix reflect the dependencies
- conducts a principle component analysis (PCA) of steps \mathbf{y}_w , sequentially in time and space
eigenvectors of the covariance matrix \mathbf{C} are the principle components / the principle axes of the mutation ellipsoid
- learns a new rotated problem representation
components are independent (only) in the new representation
- learns a new (Mahalanobis) metric
variable metric method
- approximates the inverse Hessian on quadratic functions
transformation into the sphere function
- for $\mu = 1$: conducts a natural gradient ascent on the distribution \mathcal{N} entirely independent of the given coordinate system



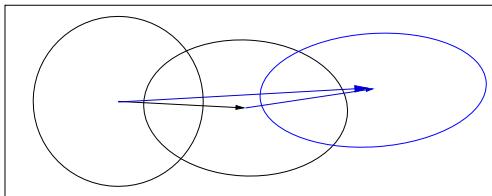
68

Cumulation

The Evolution Path

Evolution Path

Conceptually, the evolution path is the **search path** the strategy takes over a number of **generation steps**. It can be expressed as a sum of consecutive **steps** of the mean \mathbf{m} .



The recursive construction of the evolution path (cumulation):

$$\mathbf{p}_c \leftarrow \underbrace{(1 - c_e) \mathbf{p}_c}_{\text{decay factor}} + \underbrace{\sqrt{1 - (1 - c_e)^2} \sqrt{\mu_w} \mathbf{y}_w}_{\text{normalization factor}} \quad \text{input} = \frac{\mathbf{m} - \mathbf{m}_{\text{old}}}{\sigma}$$

where $\mu_w = \frac{1}{\sum w_i^2}$, $c_e \ll 1$. History information is accumulated in the evolution path.

69

“Cumulation” is a widely used technique and also known as

- *exponential smoothing* in time series, forecasting
- *exponentially weighted moving average*
- *iterate averaging* in stochastic approximation
- *momentum* in the back-propagation algorithm for ANNs
- ...

“Cumulation” conducts a *low-pass filtering*, but there is more to it...

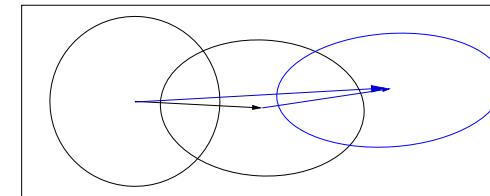
...why?

71

Cumulation

The Evolution Path

Conceptually, the evolution path is the **search path** the strategy takes over a number of **generation steps**. It can be expressed as a sum of consecutive **steps** of the mean \mathbf{m} .



The recursive construction of the evolution path (cumulation):

$$\mathbf{p}_c \leftarrow \underbrace{(1 - c_e) \mathbf{p}_c}_{\text{decay factor}} + \underbrace{\sqrt{1 - (1 - c_e)^2} \sqrt{\mu_w} \mathbf{y}_w}_{\text{normalization factor}} \quad \text{input} = \frac{\mathbf{m} - \mathbf{m}_{\text{old}}}{\sigma}$$

where $\mu_w = \frac{1}{\sum w_i^2}$, $c_e \ll 1$. History information is accumulated in the evolution path.

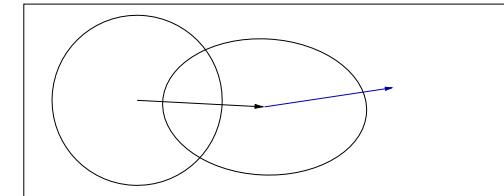
70

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mu_w \mathbf{y}_w \mathbf{y}_w^T$$

Cumulation

Utilizing the Evolution Path

We used $\mathbf{y}_w \mathbf{y}_w^T$ for updating \mathbf{C} . Because $\mathbf{y}_w \mathbf{y}_w^T = -\mathbf{y}_w (-\mathbf{y}_w)^T$ the sign of \mathbf{y}_w is lost.



The **sign information** (signifying correlation between steps) is (re-)introduced by using the *evolution path*.

$$\begin{aligned} \mathbf{p}_c &\leftarrow \underbrace{(1 - c_e) \mathbf{p}_c}_{\text{decay factor}} + \underbrace{\sqrt{1 - (1 - c_e)^2} \sqrt{\mu_w} \mathbf{y}_w}_{\text{normalization factor}} \\ \mathbf{C} &\leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \underbrace{\mathbf{p}_c \mathbf{p}_c^T}_{\text{rank-one}} \end{aligned}$$

where $\mu_w = \frac{1}{\sum w_i^2}$, $c_{\text{cov}} \ll c_e \ll 1$ such that $1/c_e$ is the “backward time horizon”.

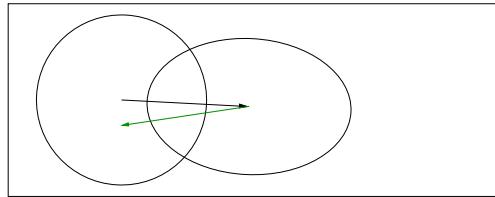
72

Cumulation

Utilizing the Evolution Path

We used $\mathbf{y}_w \mathbf{y}_w^T$ for updating \mathbf{C} . Because $\mathbf{y}_w \mathbf{y}_w^T = -\mathbf{y}_w (-\mathbf{y}_w)^T$ the sign of \mathbf{y}_w is lost.

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mu_w \mathbf{y}_w \mathbf{y}_w^T$$



The **sign information** (signifying correlation *between steps*) is (re-)introduced by using the *evolution path*.

$$\begin{aligned} \mathbf{p}_c &\leftarrow \underbrace{(1 - c_e) \mathbf{p}_c}_{\text{decay factor}} + \underbrace{\sqrt{1 - (1 - c_e)^2} \sqrt{\mu_w} \mathbf{y}_w}_{\text{normalization factor}} \\ \mathbf{C} &\leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \underbrace{\mathbf{p}_c \mathbf{p}_c^T}_{\text{rank-one}} \end{aligned}$$

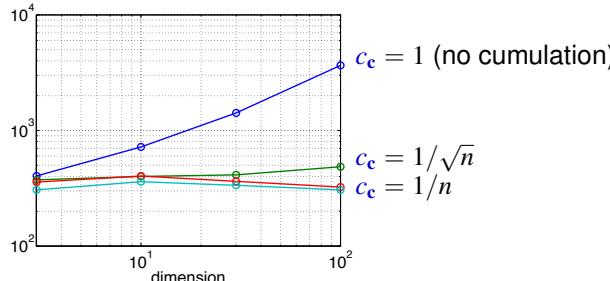
where $\mu_w = \frac{1}{\sum w_i^2}$, $c_{\text{cov}} \ll c_e \ll 1$ such that $1/c_e$ is the “backward time horizon”.

73

Using an **evolution path** for the **rank-one update** of the covariance matrix reduces the number of function evaluations to adapt to a straight ridge from about $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$.^(a)

^aHansen & Auger 2013. Principled design of continuous stochastic search: From theory to practice.

Number of f -evaluations divided by dimension on the cigar function $f(\mathbf{x}) = x_1^2 + 10^6 \sum_{i=2}^n x_i^2$



The overall model complexity is n^2 but important parts of the model can be learned in time of order n

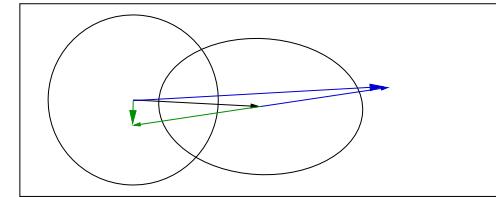
75

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mu_w \mathbf{y}_w \mathbf{y}_w^T$$

Cumulation

Utilizing the Evolution Path

We used $\mathbf{y}_w \mathbf{y}_w^T$ for updating \mathbf{C} . Because $\mathbf{y}_w \mathbf{y}_w^T = -\mathbf{y}_w (-\mathbf{y}_w)^T$ the sign of \mathbf{y}_w is lost.



The **sign information** (signifying correlation *between steps*) is (re-)introduced by using the *evolution path*.

$$\begin{aligned} \mathbf{p}_c &\leftarrow \underbrace{(1 - c_e) \mathbf{p}_c}_{\text{decay factor}} + \underbrace{\sqrt{1 - (1 - c_e)^2} \sqrt{\mu_w} \mathbf{y}_w}_{\text{normalization factor}} \\ \mathbf{C} &\leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \underbrace{\mathbf{p}_c \mathbf{p}_c^T}_{\text{rank-one}} \end{aligned}$$

where $\mu_w = \frac{1}{\sum w_i^2}$, $c_{\text{cov}} \ll c_e \ll 1$ such that $1/c_e$ is the “backward time horizon”.

74

- ➊ Problem Statement
- ➋ Evolution Strategy (ES)
- ➌ Step-Size Adaptation
 - ▶ Why Step-Size Control
 - ▶ Path Length Control (CSA)
 - ▶ Limitations of CSA and its Alternatives
- ➍ Covariance Matrix Adaptation (CMA)
 - ▶ Rank-One Update and Cumulation
 - ▶ Rank- μ Update
 - ▶ Active Covariance Update
- ➎ Design Principle
 - ▶ Theoretical Foundations
 - ▶ Variants for Large Scale Problems
- ➏ Summary and Final Remarks

76

Rank- μ Update

$$\begin{aligned} \mathbf{x}_i &= \mathbf{m} + \sigma \mathbf{y}_i, & \mathbf{y}_i &\sim \mathcal{N}_i(\mathbf{0}, \mathbf{C}), \\ \mathbf{m} &\leftarrow \mathbf{m} + \sigma \mathbf{y}_w & \mathbf{y}_w &= \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \end{aligned}$$

The rank- μ update extends the update rule for **large population sizes** λ using $\mu > 1$ vectors to update \mathbf{C} at each generation step.

The weighted empirical covariance matrix

$$\mathbf{C}_\mu = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

computes a weighted mean of the outer products of the best μ steps and has rank $\min(\mu, n)$ with probability one.

with $\mu = \lambda$ weights can be negative ¹⁰

The rank- μ update then reads

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mathbf{C}_\mu$$

where $c_{\text{cov}} \approx \mu_w/n^2$ and $c_{\text{cov}} \leq 1$.

¹⁰ Jastrebski and Arnold (2006). Improving evolution strategies through active covariance matrix adaptation. CEC. 77

Rank- μ Update

$$\begin{aligned} \mathbf{x}_i &= \mathbf{m} + \sigma \mathbf{y}_i, & \mathbf{y}_i &\sim \mathcal{N}_i(\mathbf{0}, \mathbf{C}), \\ \mathbf{m} &\leftarrow \mathbf{m} + \sigma \mathbf{y}_w & \mathbf{y}_w &= \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \end{aligned}$$

The rank- μ update extends the update rule for **large population sizes** λ using $\mu > 1$ vectors to update \mathbf{C} at each generation step.

The weighted empirical covariance matrix

$$\mathbf{C}_\mu = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

computes a weighted mean of the outer products of the best μ steps and has rank $\min(\mu, n)$ with probability one.

with $\mu = \lambda$ weights can be negative ¹⁰

The rank- μ update then reads

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mathbf{C}_\mu$$

where $c_{\text{cov}} \approx \mu_w/n^2$ and $c_{\text{cov}} \leq 1$.

¹⁰ Jastrebski and Arnold (2006). Improving evolution strategies through active covariance matrix adaptation. CEC. 79

Rank- μ Update

$$\begin{aligned} \mathbf{x}_i &= \mathbf{m} + \sigma \mathbf{y}_i, & \mathbf{y}_i &\sim \mathcal{N}_i(\mathbf{0}, \mathbf{C}), \\ \mathbf{m} &\leftarrow \mathbf{m} + \sigma \mathbf{y}_w & \mathbf{y}_w &= \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \end{aligned}$$

The rank- μ update extends the update rule for **large population sizes** λ using $\mu > 1$ vectors to update \mathbf{C} at each generation step.

The weighted empirical covariance matrix

$$\mathbf{C}_\mu = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

computes a weighted mean of the outer products of the best μ steps and has rank $\min(\mu, n)$ with probability one.

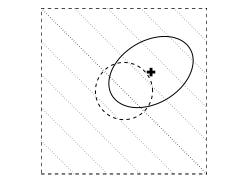
with $\mu = \lambda$ weights can be negative ¹⁰

The rank- μ update then reads

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mathbf{C}_\mu$$

where $c_{\text{cov}} \approx \mu_w/n^2$ and $c_{\text{cov}} \leq 1$.

¹⁰ Jastrebski and Arnold (2006). Improving evolution strategies through active covariance matrix adaptation. CEC. 78



$$\mathbf{x}_i = \mathbf{m} + \sigma \mathbf{y}_i, \quad \mathbf{y}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$$

$$\mathbf{C}_\mu = \frac{1}{\mu} \sum_{i=1}^{\mu} \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

$$\mathbf{m}_{\text{new}} \leftarrow \mathbf{m} + \frac{1}{\mu} \sum_{i=1}^{\mu} \mathbf{y}_{i:\lambda}$$

new distribution

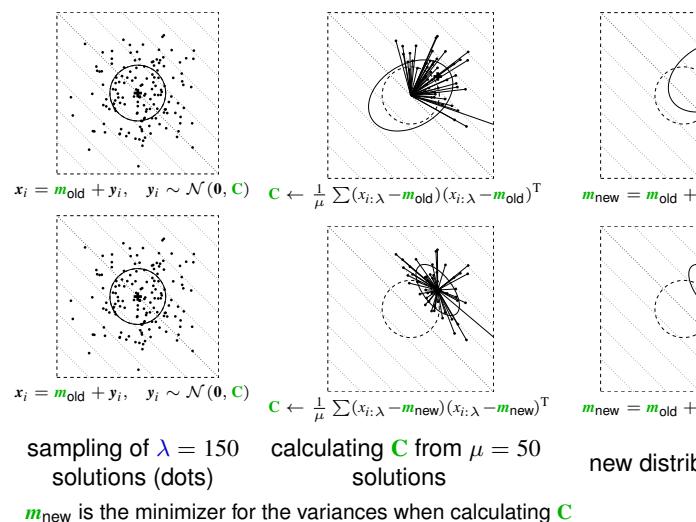
sampling of $\lambda = 150$ solutions where

$$\mathbf{C} = \mathbf{I} \text{ and } \sigma = 1$$

calculating \mathbf{C} where $\mu = 50$,

$$w_1 = \dots = w_\mu = \frac{1}{\mu}, \text{ and } c_{\text{cov}} = 1$$

Rank- μ CMA versus Estimation of Multivariate Normal Algorithm EMNA_{global}¹¹



¹¹ Hansen, N. (2006). The CMA Evolution Strategy: A Comparing Review. In J.A. Lozano, P. Larranga, I. Inza and E. Bengoechea (Eds.). Towards a new evolutionary computation. Advances in estimation of distribution algorithms. pp. 75-102
81

The rank- μ update

- increases the possible learning rate in large populations
roughly from $2/n^2$ to μ_w/n^2
- can reduce the number of necessary generations roughly from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ ⁽¹²⁾
given $\mu_w \propto \lambda \propto n$

Therefore the rank- μ update is the primary mechanism whenever a large population size is used

say $\lambda \geq 3n + 10$

The rank-one update

- uses the evolution path and reduces the number of necessary function evaluations to learn straight ridges from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$.

Rank-one update and rank- μ update can be combined

...all equations

¹²Hansen, Müller, and Koumoutsakos 2003. Reducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation*, 11(1), pp. 1-18
83

The rank- μ update

- increases the possible learning rate in large populations
roughly from $2/n^2$ to μ_w/n^2
- can reduce the number of necessary generations roughly from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ ⁽¹²⁾
given $\mu_w \propto \lambda \propto n$

Therefore the rank- μ update is the primary mechanism whenever a large population size is used

say $\lambda \geq 3n + 10$

The rank-one update

- uses the evolution path and reduces the number of necessary function evaluations to learn straight ridges from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$.

Rank-one update and rank- μ update can be combined

...all equations

¹²Hansen, Müller, and Koumoutsakos 2003. Reducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation*, 11(1), pp. 1-18
82

The rank- μ update

- increases the possible learning rate in large populations
roughly from $2/n^2$ to μ_w/n^2
- can reduce the number of necessary generations roughly from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ ⁽¹²⁾
given $\mu_w \propto \lambda \propto n$

Therefore the rank- μ update is the primary mechanism whenever a large population size is used

say $\lambda \geq 3n + 10$

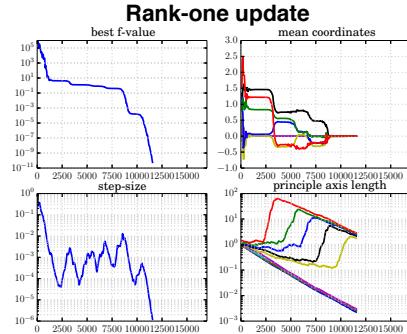
The rank-one update

- uses the evolution path and reduces the number of necessary function evaluations to learn straight ridges from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$.

Rank-one update and rank- μ update can be combined

...all equations

¹²Hansen, Müller, and Koumoutsakos 2003. Reducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation*, 11(1), pp. 1-18
84



$$f_{\text{TwoAxes}}(x) = \sum_{i=1}^5 x_i^2 + 10^6 \sum_{i=6}^{10} x_i^2$$

$\lambda = 10$ (default for $N = 10$)

85

1 Problem Statement

2 Evolution Strategy (ES)

3 Step-Size Adaptation

- ▶ Why Step-Size Control
- ▶ Path Length Control (CSA)
- ▶ Limitations of CSA and its Alternatives

4 Covariance Matrix Adaptation (CMA)

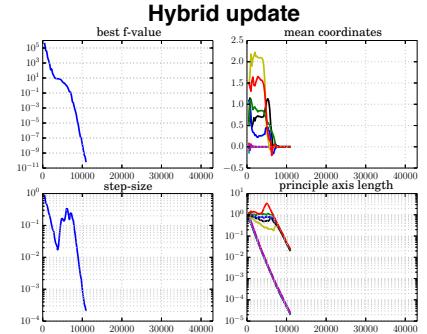
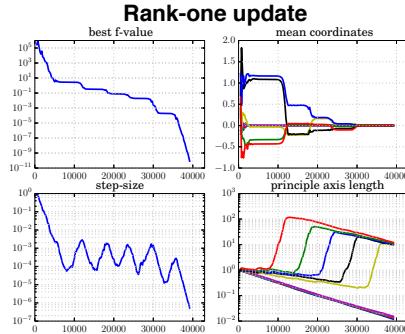
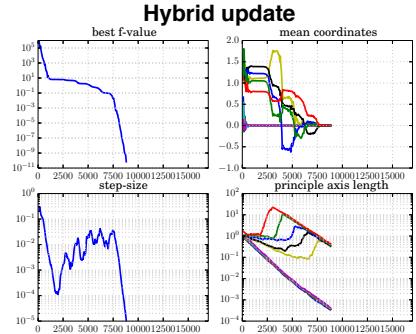
- ▶ Rank-One Update and Cumulation
- ▶ Rank- μ Update
- ▶ Active Covariance Update

5 Design Principle

- ▶ Theoretical Foundations
- ▶ Variants for Large Scale Problems

6 Summary and Final Remarks

87



$$f_{\text{TwoAxes}}(x) = \sum_{i=1}^5 x_i^2 + 10^6 \sum_{i=6}^{10} x_i^2$$

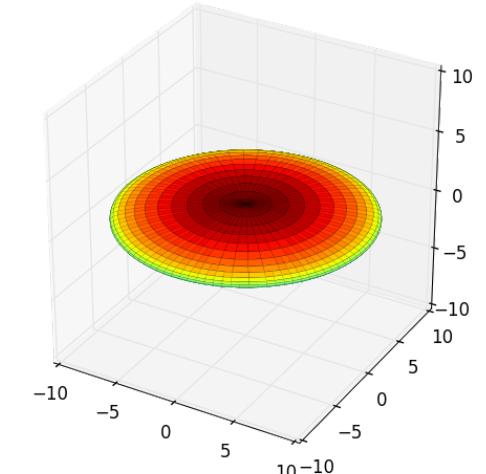
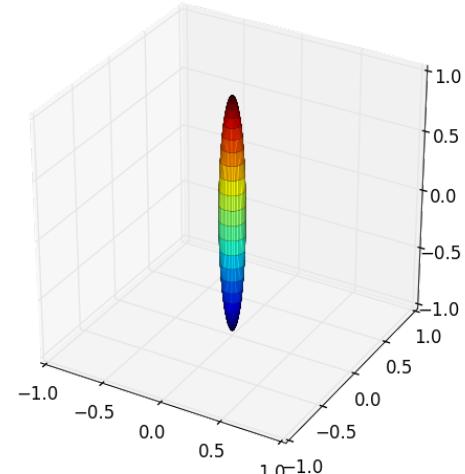
$\lambda = 50$

86

Different Types of Ill-Conditioning

(Axes Ratio = 10)

Cigar Type:
1 long axis = n-1 short axes



Discus Type:
1 short axis = n-1 long axes

Learning a Short/Long Axis

Rank-one and Rank- μ Covariance Matrix Update

$$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \underbrace{\sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T}_{\text{nonnegative definite, i.e., increase covariances}}$$

- increases the variance in the directions of \mathbf{p}_c and $\mathbf{y}_{i:\lambda}$; cumulation is very effective for this purpose
- decrease the variance only by multiplying the factor $(1 - c_1 - c_\mu)$

89

Learning a Short/Long Axis

Rank-one and Rank- μ Covariance Matrix Update

$$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \underbrace{\sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T}_{\text{nonnegative definite, i.e., increase covariances}}$$

- increases the variance in the directions of \mathbf{p}_c and $\mathbf{y}_{i:\lambda}$; cumulation is very effective for this purpose
- decrease the variance only by multiplying the factor $(1 - c_1 - c_\mu)$

90

Learning a Short/Long Axis

Rank-one and Rank- μ Covariance Matrix Update

$$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \underbrace{\sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T}_{\text{nonnegative definite, i.e., increase covariances}}$$

- increases the variance in the directions of \mathbf{p}_c and $\mathbf{y}_{i:\lambda}$; cumulation is very effective for this purpose
- decrease the variance only by multiplying the factor $(1 - c_1 - c_\mu)$

91

Learning a Short/Long Axis

Rank-one and Rank- μ Covariance Matrix Update

$$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \underbrace{\sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T}_{\text{nonnegative definite, i.e., increase covariances}}$$

92

Active Update

utilize negative weights [Jastrebski and Arnold, 2006]

Rank-one and Rank- μ Active Covariance Matrix Update

$$\mathbf{C} \leftarrow (1 - c_1 - c_\mu + c_\mu^-) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^{\mu} \underbrace{\mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T}_{\text{nonnegative definite}} - c_\mu^- \sum_{i=\lambda-\mu}^{\lambda} |\mathbf{w}_i| \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

where

- $-|\mathbf{w}_i| < 0$ (for $i \geq \lambda - \mu$): negative weight assigned to $\mathbf{y}_{i:\lambda}$, $\sum_{i=\lambda-\mu}^{\lambda} |\mathbf{w}_i| = 1$.
- $c_\mu^- > 0$: learning rate for the active update

[Jastrebski and Arnold, 2006] Jastrebski, G. and Arnold, D. V. (2006). Improving Evolution Strategies through Active Covariance Matrix Adaptation. In 2006 IEEE Congress on Evolutionary Computation, pages 9719–9726.

93

Single Update on Discus Function ($n = 3$)

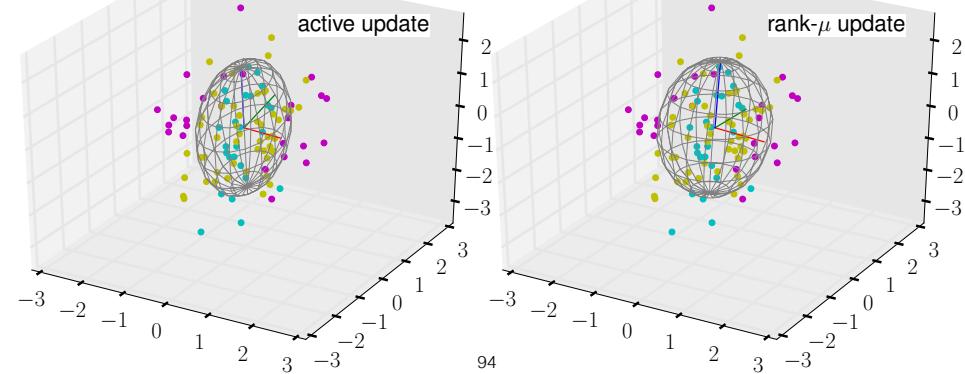
$$\mathbf{C} \leftarrow (1 - c_1 - c_\mu^- + c_\mu^-) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T - c_\mu^- \sum_{i=\lambda-\mu}^{\lambda} |\mathbf{w}_i| \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

- $\lambda = 100$, $\mu = \lambda/4$, $\mathbf{w}_i = 1/\mu$

- $c_\mu = 0.3$, $c_\mu^- = 0$ for rank- μ update and $c_\mu^- = 0.2$ for active update

$u_1 = [0.9998, 0.0168, -0.0080]$, $d_1 = 6.35e-01$
 $u_2 = [0.0160, 0.9948, 0.1002]$, $d_2 = 9.81e-01$
 $u_3 = [0.0097, -0.1001, 0.9949]$, $d_3 = 1.13e+00$

$u_1 = [1.0000, -0.0010, -0.0004]$, $d_1 = 8.44e-01$
 $u_2 = [0.0010, 0.9887, -0.1501]$, $d_2 = 9.67e-01$
 $u_3 = [0.0005, 0.1501, 0.9887]$, $d_3 = 1.10e+00$



Summary

Active Covariance Matrix Adaptation + Cumulation

$$\mathbf{C} \leftarrow (1 - c_1 - c_\mu + c_\mu^-) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T - c_\mu^- \sum_{i=\lambda-\mu}^{\lambda} |\mathbf{w}_i| \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

- $-|\mathbf{w}_i| < 0$ (for $i \geq \lambda - \mu$): negative weight assigned to $\mathbf{y}_{i:\lambda}$, $\sum_{i=\lambda-\mu}^{\lambda} |\mathbf{w}_i| = 1$.
- $c_\mu^- > 0$: learning rate for the active update

These components compensate each other

- cumulation: excels to learn a long axis, but inefficient for a large λ
- rank- μ update: efficient for a large λ
- active update: effective to learn short axes

An important yet solvable issue of active update

- The positive definiteness of \mathbf{C} will be violated if c_μ^- is not small enough
- The positive definiteness can be guaranteed w.p.1 by controlling $c_\mu^- \mathbf{w}_i$

95

556

Summary

Active Covariance Matrix Adaptation + Cumulation

$$\mathbf{C} \leftarrow (1 - c_1 - c_\mu + c_\mu^-) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T - c_\mu^- \sum_{i=\lambda-\mu}^{\lambda} |\mathbf{w}_i| \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

- $-|\mathbf{w}_i| < 0$ (for $i \geq \lambda - \mu$): negative weight assigned to $\mathbf{y}_{i:\lambda}$, $\sum_{i=\lambda-\mu}^{\lambda} |\mathbf{w}_i| = 1$.
- $c_\mu^- > 0$: learning rate for the active update

These components compensate each other

- cumulation: excels to learn a long axis, but inefficient for a large λ
- rank- μ update: efficient for a large λ
- active update: effective to learn short axes

An important yet solvable issue of active update

- The positive definiteness of \mathbf{C} will be violated if c_μ^- is not small enough
- The positive definiteness can be guaranteed w.p.1 by controlling $c_\mu^- \mathbf{w}_i$

96

Summary

Active Covariance Matrix Adaptation + Cumulation

$$\mathbf{C} \leftarrow (1 - c_1 - c_\mu + c_\mu^-) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T - c_\mu^- \sum_{i=\lambda-\mu}^{\lambda} |\mathbf{w}_i| \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

- $-|\mathbf{w}_i| < 0$ (for $i \geq \lambda - \mu$): negative weight assigned to $\mathbf{y}_{i:\lambda}$, $\sum_{i=\lambda-\mu}^{\lambda} |\mathbf{w}_i| = 1$.
- $c_\mu^- > 0$: learning rate for the active update

These components compensate each other

- cumulation: excels to learn a long axis, but inefficient for a large λ
- rank- μ update: efficient for a large λ
- active update: effective to learn short axes

An important yet solvable issue of active update

- The positive definiteness of \mathbf{C} will be violated if c_μ^- is not small enough
- The positive definiteness can be guaranteed w.p.1 by controlling $c_\mu^- \mathbf{w}_i$

97

Strategy Internal Parameters

- related to selection and recombination
 - ▶ λ , offspring number, new solutions sampled, population size
 - ▶ μ , parent number, solutions involved in updates of \mathbf{m} , \mathbf{C} , and σ
 - ▶ $w_{i=1,\dots,\mu}$, recombination weights
- related to \mathbf{C} -update
 - ▶ c_e , decay rate for the evolution path
 - ▶ c_1 , learning rate for rank-one update of \mathbf{C}
 - ▶ c_μ , learning rate for rank- μ update of \mathbf{C}
- related to σ -update
 - ▶ c_σ , decay rate of the evolution path
 - ▶ d_σ , damping for σ -change

Parameters were identified in carefully chosen experimental set ups. Parameters do not in the

first place depend on the objective function and are not meant to be in the users choice.

Only(?) the population size λ (and the initial σ) might be reasonably varied in a wide range, depending on the objective function

Useful: restarts with increasing population size (IPOP)

99

Summary of Equations

The Covariance Matrix Adaptation Evolution Strategy

Input: $\mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, λ (problem dependent)

Initialize: $\mathbf{C} = \mathbf{I}$, and $\mathbf{p}_c = \mathbf{0}$, $\mathbf{p}_\sigma = \mathbf{0}$,

Set: $c_e \approx 4/n$, $c_\sigma \approx 4/n$, $c_1 \approx 2/n^2$, $c_\mu \approx \mu_w/n^2$, $c_1 + c_\mu \leq 1$, $d_\sigma \approx 1 + \sqrt{\frac{\mu_w}{n}}$, and $w_{i=1,\dots,\lambda}$ such that $\mu_w = \frac{1}{\sum_{i=1}^{\mu} w_i^2} \approx 0.3 \lambda$

While not terminate

$$\mathbf{x}_i = \mathbf{m} + \sigma \mathbf{y}_i, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C}), \quad \text{for } i = 1, \dots, \lambda$$

$$\mathbf{m} \leftarrow \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{x}_{i:\lambda} = \mathbf{m} + \sigma \mathbf{y}_w \quad \text{where } \mathbf{y}_w = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda}$$

$$\mathbf{p}_c \leftarrow (1 - c_e) \mathbf{p}_c + \mathbf{1}_{\{\|\mathbf{p}_c\| < 1.5\sqrt{n}\}} \sqrt{1 - (1 - c_e)^2} \sqrt{\mu_w} \mathbf{y}_w$$

$$\mathbf{p}_\sigma \leftarrow (1 - c_\sigma) \mathbf{p}_\sigma + \sqrt{1 - (1 - c_\sigma)^2} \sqrt{\mu_w} \mathbf{C}^{-\frac{1}{2}} \mathbf{y}_w$$

$$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

$$\sigma \leftarrow \sigma \times \exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E}[\|\mathcal{N}(0, \mathbf{I})\|]} - 1 \right) \right)$$

sampling

update mean

cumulation for \mathbf{C} cumulation for σ update \mathbf{C} update of σ

Not covered on this slide: termination, restarts, useful output, boundaries and encoding

98



- ① Problem Statement
- ② Evolution Strategy (ES)
- ③ Step-Size Adaptation
 - ▶ Why Step-Size Control
 - ▶ Path Length Control (CSA)
 - ▶ Limitations of CSA and its Alternatives
- ④ Covariance Matrix Adaptation (CMA)
 - ▶ Rank-One Update and Cumulation
 - ▶ Rank- μ Update
 - ▶ Active Covariance Update
- ⑤ Design Principle
 - ▶ Theoretical Foundations
 - ▶ Variants for Large Scale Problems
- ⑥ Summary and Final Remarks

557

100



Natural Gradient Descend

- Consider $\arg \min_{\theta} E(f(x)|\theta)$ under the sampling distribution $x \sim p(\cdot|\theta)$
we could improve $E(f(x)|\theta)$ by following the gradient $\nabla_{\theta} E(f(x)|\theta)$:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} E(f(x)|\theta), \quad \eta > 0$$

∇_{θ} depends on the parameterization of the distribution, therefore

- Consider the **natural gradient** of the expected transformed fitness

$$\begin{aligned}\tilde{\nabla}_{\theta} E(w \circ P_f(f(x))|\theta) &= F_{\theta}^{-1} \nabla_{\theta} E(w \circ P_f(f(x))|\theta) \\ &= E(w \circ P_f(f(x)) F_{\theta}^{-1} \nabla_{\theta} \ln p(x|\theta))\end{aligned}$$

using the Fisher information matrix $F_{\theta} = \left(\left(E \frac{\partial^2 \log p(x|\theta)}{\partial \theta_i \partial \theta_j} \right) \right)_{ij}$ of the density p .

The natural gradient is **invariant under re-parameterization** of the distribution.

- A Monte-Carlo approximation reads

$$\tilde{\nabla}_{\theta} \hat{E}(\hat{w}(f(x))|\theta) = \sum_{i=1}^{\lambda} w_i F_{\theta}^{-1} \nabla_{\theta} \ln p(x_{i:\lambda}|\theta), \quad w_i = \hat{w}(f(x_{i:\lambda})|\theta)$$

101



Natural Gradient Descend

- Consider $\arg \min_{\theta} E(f(x)|\theta)$ under the sampling distribution $x \sim p(\cdot|\theta)$
we could improve $E(f(x)|\theta)$ by following the gradient $\nabla_{\theta} E(f(x)|\theta)$:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} E(f(x)|\theta), \quad \eta > 0$$

∇_{θ} depends on the parameterization of the distribution, therefore

- Consider the **natural gradient** of the expected transformed fitness

$$\begin{aligned}\tilde{\nabla}_{\theta} E(w \circ P_f(f(x))|\theta) &= F_{\theta}^{-1} \nabla_{\theta} E(w \circ P_f(f(x))|\theta) \\ &= E(w \circ P_f(f(x)) F_{\theta}^{-1} \nabla_{\theta} \ln p(x|\theta))\end{aligned}$$

using the Fisher information matrix $F_{\theta} = \left(\left(E \frac{\partial^2 \log p(x|\theta)}{\partial \theta_i \partial \theta_j} \right) \right)_{ij}$ of the density p .

The natural gradient is **invariant under re-parameterization** of the distribution.

- A Monte-Carlo approximation reads

$$\tilde{\nabla}_{\theta} \hat{E}(\hat{w}(f(x))|\theta) = \sum_{i=1}^{\lambda} w_i F_{\theta}^{-1} \nabla_{\theta} \ln p(x_{i:\lambda}|\theta), \quad w_i = \hat{w}(f(x_{i:\lambda})|\theta)$$

103



Natural Gradient Descend

- Consider $\arg \min_{\theta} E(f(x)|\theta)$ under the sampling distribution $x \sim p(\cdot|\theta)$
we could improve $E(f(x)|\theta)$ by following the gradient $\nabla_{\theta} E(f(x)|\theta)$:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} E(f(x)|\theta), \quad \eta > 0$$

∇_{θ} depends on the parameterization of the distribution, therefore

- Consider the **natural gradient** of the expected transformed fitness

$$\begin{aligned}\tilde{\nabla}_{\theta} E(w \circ P_f(f(x))|\theta) &= F_{\theta}^{-1} \nabla_{\theta} E(w \circ P_f(f(x))|\theta) \\ &= E(w \circ P_f(f(x)) F_{\theta}^{-1} \nabla_{\theta} \ln p(x|\theta))\end{aligned}$$

using the Fisher information matrix $F_{\theta} = \left(\left(E \frac{\partial^2 \log p(x|\theta)}{\partial \theta_i \partial \theta_j} \right) \right)_{ij}$ of the density p .

The natural gradient is **invariant under re-parameterization** of the distribution.

- A Monte-Carlo approximation reads

$$\tilde{\nabla}_{\theta} \hat{E}(\hat{w}(f(x))|\theta) = \sum_{i=1}^{\lambda} w_i F_{\theta}^{-1} \nabla_{\theta} \ln p(x_{i:\lambda}|\theta), \quad w_i = \hat{w}(f(x_{i:\lambda})|\theta)$$

102



Natural Gradient Descend

- Consider $\arg \min_{\theta} E(f(x)|\theta)$ under the sampling distribution $x \sim p(\cdot|\theta)$
we could improve $E(f(x)|\theta)$ by following the gradient $\nabla_{\theta} E(f(x)|\theta)$:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} E(f(x)|\theta), \quad \eta > 0$$

∇_{θ} depends on the parameterization of the distribution, therefore

- Consider the **natural gradient** of the expected transformed fitness

$$\begin{aligned}\tilde{\nabla}_{\theta} E(w \circ P_f(f(x))|\theta) &= F_{\theta}^{-1} \nabla_{\theta} E(w \circ P_f(f(x))|\theta) \\ &= E(w \circ P_f(f(x)) F_{\theta}^{-1} \nabla_{\theta} \ln p(x|\theta))\end{aligned}$$

using the Fisher information matrix $F_{\theta} = \left(\left(E \frac{\partial^2 \log p(x|\theta)}{\partial \theta_i \partial \theta_j} \right) \right)_{ij}$ of the density p .

The natural gradient is **invariant under re-parameterization** of the distribution.

- A Monte-Carlo approximation reads

$$\tilde{\nabla}_{\theta} \hat{E}(\hat{w}(f(x))|\theta) = \sum_{i=1}^{\lambda} w_i F_{\theta}^{-1} \nabla_{\theta} \ln p(x_{i:\lambda}|\theta), \quad w_i = \hat{w}(f(x_{i:\lambda})|\theta)$$

104



Natural Gradient Descend

- Consider $\arg \min_{\theta} E(f(\mathbf{x})|\theta)$ under the sampling distribution $\mathbf{x} \sim p(\cdot|\theta)$
we could improve $E(f(\mathbf{x})|\theta)$ by following the gradient $\nabla_{\theta} E(f(\mathbf{x})|\theta)$:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} E(f(\mathbf{x})|\theta), \quad \eta > 0$$

∇_{θ} depends on the parameterization of the distribution, therefore

- Consider the **natural gradient** of the expected transformed fitness

$$\begin{aligned}\tilde{\nabla}_{\theta} E(w \circ P_f(f(\mathbf{x}))|\theta) &= F_{\theta}^{-1} \nabla_{\theta} E(w \circ P_f(f(\mathbf{x}))|\theta) \\ &= E(w \circ P_f(f(\mathbf{x}))) F_{\theta}^{-1} \nabla_{\theta} \ln p(\mathbf{x}|\theta))\end{aligned}$$

using the Fisher information matrix $F_{\theta} = \left(\left(E \frac{\partial^2 \log p(\mathbf{x}|\theta)}{\partial \theta_i \partial \theta_j} \right) \right)_{ij}$ of the density p .

The natural gradient is **invariant under re-parameterization** of the distribution.

- A **Monte-Carlo approximation** reads

$$\tilde{\nabla}_{\theta} \widehat{E}(\widehat{w}(f(\mathbf{x}))|\theta) = \sum_{i=1}^{\lambda} w_i F_{\theta}^{-1} \nabla_{\theta} \ln p(\mathbf{x}_{i:\lambda}|\theta), \quad w_i = \widehat{w}(f(\mathbf{x}_{i:\lambda})|\theta)$$

105

Maximum Likelihood Update

The new distribution mean \mathbf{m} maximizes the log-likelihood

$$\mathbf{m}_{\text{new}} = \arg \max_{\mathbf{m}} \sum_{i=1}^{\mu} w_i \log p_{\mathcal{N}}(\mathbf{x}_{i:\lambda}|\mathbf{m})$$

independently of the given covariance matrix

$$\log p_{\mathcal{N}}(\mathbf{x}|\mathbf{m}, \mathbf{C}) = -\frac{1}{2} \log \det(2\pi\mathbf{C}) - \frac{1}{2} (\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m})$$

$p_{\mathcal{N}}$ is the density of the multi-variate normal distribution

107

CMA-ES = Natural Evolution Strategy + Cumulation

Natural gradient descend using the MC approximation and the normal distribution

- Rewriting the update of the distribution mean

$$\mathbf{m}_{\text{new}} \leftarrow \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{x}_{i:\lambda} = \mathbf{m} + \underbrace{\sum_{i=1}^{\mu} \mathbf{w}_i (\mathbf{x}_{i:\lambda} - \mathbf{m})}_{\text{natural gradient for mean } \frac{\partial}{\partial \mathbf{m}} \widehat{E}(w \circ P_f(f(\mathbf{x})))|\mathbf{m}, \mathbf{C})}$$

- Rewriting the update of the covariance matrix¹³

$$\begin{aligned}\mathbf{C}_{\text{new}} &\leftarrow \mathbf{C} + \underbrace{\mathbf{c}_1 (\mathbf{p} \mathbf{c}^T - \mathbf{C})}_{\text{rank one}} \\ &+ \underbrace{\frac{\mathbf{c}_\mu}{\sigma^2} \sum_{i=1}^{\mu} \mathbf{w}_i \left((\mathbf{x}_{i:\lambda} - \mathbf{m}) (\mathbf{x}_{i:\lambda} - \mathbf{m})^T - \sigma^2 \mathbf{C} \right)}_{\text{rank-}\mu} \\ &\text{natural gradient for covariance matrix } \frac{\partial}{\partial \mathbf{C}} \widehat{E}(w \circ P_f(f(\mathbf{x})))|\mathbf{m}, \mathbf{C})\end{aligned}$$

¹³ Akimoto et.al. (2010): Bidirectional Relation between CMA Evolution Strategies and Natural Evolution Strategies, PPSN XI, 106

Maximum Likelihood Update

The new distribution mean \mathbf{m} maximizes the log-likelihood

$$\mathbf{m}_{\text{new}} = \arg \max_{\mathbf{m}} \sum_{i=1}^{\mu} w_i \log p_{\mathcal{N}}(\mathbf{x}_{i:\lambda}|\mathbf{m})$$

independently of the given covariance matrix

The rank- μ update matrix \mathbf{C}_{μ} maximizes the log-likelihood

$$\mathbf{C}_{\mu} = \arg \max_{\mathbf{C}} \sum_{i=1}^{\mu} w_i \log p_{\mathcal{N}} \left(\frac{\mathbf{x}_{i:\lambda} - \mathbf{m}_{\text{old}}}{\sigma} \middle| \mathbf{m}_{\text{old}}, \mathbf{C} \right)$$

$$\log p_{\mathcal{N}}(\mathbf{x}|\mathbf{m}, \mathbf{C}) = -\frac{1}{2} \log \det(2\pi\mathbf{C}) - \frac{1}{2} (\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m})$$

$p_{\mathcal{N}}$ is the density of the multi-variate normal distribution

108

CMA-ES as a Natural Gradient

$\mathcal{O}(\lambda N^2)$ floating point (FP) multiplication + $\mathcal{O}(N^2 + \lambda N)$ FP memory. ($\lambda \in o(N)$)

1. $\sqrt{\mathbf{C}^{(t)}} = \text{MATRIXSQRT}(\mathbf{C}^{(t)})$ (perform every $\mathcal{O}(N/\lambda)$ iter.)
2. $z_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ for $i = 1, \dots, \lambda$
3. $x_i = \mathbf{m}^{(t)} + \sigma^{(t)} \sqrt{\mathbf{C}^{(t)}} z_i$
4. $(x_{i:\lambda})_{i=1,\dots,\lambda} = \text{SORTW.R.T.}^f((x_i)_{i=1,\dots,\lambda})$
5. $\mathbf{p}_c^{(t+1)} = (1 - c_c)\mathbf{p}_c^{(t)} + \sqrt{c_c(2 - c_c)/(\sum w_i^2)} \sum_{i=1}^\lambda w_i (x_{i:\lambda} - \mathbf{m}^{(t)}) / \sigma^{(t)}$
6. $\mathbf{p}_\sigma^{(t+1)} = (1 - c_\sigma)\mathbf{p}_\sigma^{(t)} + \sqrt{c_\sigma(2 - c_\sigma)/(\sum w_i^2)} \sum_{i=1}^\lambda w_i z_{i:\lambda}$
7. $\sigma^{(t+1)} = \sigma^{(t)} \exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma^{(t+1)}\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)$
8. $\mathbf{m}^{(t+1)} = \mathbf{m}^{(t)} + \sum_{i=1}^\lambda w_i \tilde{\nabla} J(x_{i:\lambda})$
9. $\mathbf{C}^{(t+1)} = \mathbf{C}^{(t)} + c_\mu \sum_{i=1}^\lambda w_i \tilde{\nabla} J(x_{i:\lambda}) + c_1 \tilde{\nabla} J(\mathbf{m}^{(t)} + \sigma^{(t)} \mathbf{p}_c^{(t+1)})$

$\mathcal{O}(\lambda N^2)$

$\mathcal{O}(\lambda N)$

$\mathcal{O}(\lambda N^2)$

109

Bottleneck of the CMA-ES

- ➊ $\mathcal{O}(N^2)$ Time and Space Complexities
 - ▶ to store and update $\mathbf{C} \in \mathbb{R}^{N \times N}$
 - ▶ to compute the eigen decomposition of \mathbf{C}
- ➋ $\mathcal{O}(1/N^2)$ Learning Rates for \mathbf{C} -Update
 - ▶ $c_\mu \approx \mu_w / N^2$
 - ▶ $c_1 \approx 2 / N^2$

⇒ Design a variant of CMA-ES for high dimensional search space

110

Variants with Restricted Covariance Matrix

CMA-ES Variants with Restricted Covariance Matrices

- Sep-CMA [Ros and Hansen, 2008]
 - ▶ $\mathbf{C} = \mathbf{D}$. \mathbf{D} : Diagonal
- VD-CMA [Akimoto et al., 2014]
 - ▶ $\mathbf{C} = \mathbf{D}(\mathbf{I} + \mathbf{v}\mathbf{v}^T)\mathbf{D}$. \mathbf{D} : Diagonal, $\mathbf{v} \in \mathbb{R}^N$.
- LM-CMA [Loshchilov, 2014]
 - ▶ $\mathbf{C} = \mathbf{I} + \sum_{i=1}^k \mathbf{v}_i \mathbf{v}_i^T$. $\mathbf{v}_i \in \mathbb{R}^N$.
- VkD-CMA [Akimoto and Hansen, 2016]
 - ▶ $\mathbf{C} = \mathbf{D}(\mathbf{I} + \sum_{i=1}^k \mathbf{v}_i \mathbf{v}_i^T)\mathbf{D}$. $\mathbf{v}_i \in \mathbb{R}^N$.

[Ros and Hansen, 2008] Ros, R. and Hansen, N. (2008). A simple modification in CMA-ES achieving linear time and space complexity. In Parallel Problem Solving from Nature - PPSN X, pages 296–305. Springer.

[Akimoto et al., 2014] Akimoto, Y., Auger, A., and Hansen, N. (2014). Comparison-based natural gradient optimization in high dimension. In Proceedings of Genetic and Evolutionary Computation Conference, pages 373–380, Vancouver, BC, Canada.

[Loshchilov, 2014] Loshchilov, I. (2014). A computationally efficient limited memory cma-es for large scale optimization. In Proceedings of Genetic and Evolutionary Computation Conference, pages 397–404.

[Akimoto and Hansen, 2016] Akimoto, Y. and Hansen, N. (2016). Projection-based restricted covariance matrix adaptation for high dimension. In Genetic and Evolutionary Computation Conference, GECCO 2016, Denver, Colorado, USA, July 20-24, 2016, page (accepted). ACM.

111

560

112

CMA-ES as a Natural Gradient

1. $\sqrt{\mathbf{C}^{(t)}} = \text{MATRIXSQRT}(\mathbf{C}^{(t)})$ (perform every $\mathcal{O}(N/\lambda)$ iter.)
2. $z_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ for $i = 1, \dots, \lambda$
3. $x_i = \mathbf{m}^{(t)} + \sigma^{(t)} \sqrt{\mathbf{C}^{(t)}} z_i$
4. $(x_{i:\lambda})_{i=1,\dots,\lambda} = \text{SORTW.R.T.}^f((x_i)_{i=1,\dots,\lambda})$
5. $\mathbf{p}_c^{(t+1)} = (1 - c_c)\mathbf{p}_c^{(t)} + \sqrt{c_c(2 - c_c)/(\sum w_i^2)} \sum_{i=1}^\lambda w_i (x_{i:\lambda} - \mathbf{m}^{(t)}) / \sigma^{(t)}$
6. $\mathbf{p}_\sigma^{(t+1)} = (1 - c_\sigma)\mathbf{p}_\sigma^{(t)} + \sqrt{c_\sigma(2 - c_\sigma)/(\sum w_i^2)} \sum_{i=1}^\lambda w_i z_{i:\lambda}$
7. $\sigma^{(t+1)} = \sigma^{(t)} \exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma^{(t+1)}\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)$
8. $\mathbf{m}^{(t+1)} = \mathbf{m}^{(t)} + \sum_{i=1}^\lambda w_i \tilde{\nabla} J(x_{i:\lambda})$
9. $\mathbf{C}^{(t+1)} = \mathbf{C}^{(t)} + c_\mu \sum_{i=1}^\lambda w_i \tilde{\nabla} J(x_{i:\lambda}) + c_1 \tilde{\nabla} J(\mathbf{m}^{(t)} + \sigma^{(t)} \mathbf{p}_c^{(t+1)})$

CMA-ES as a Natural Gradient

$\mathcal{O}(\lambda N^2)$ floating point (FP) multiplication + $\mathcal{O}(N^2 + \lambda N)$ FP memory. ($\lambda \in o(N)$)

$$1. \sqrt{\mathbf{C}^{(t)}} = \text{MATRIXSQRT}(\mathbf{C}^{(t)}) \quad (\text{perform every } \mathcal{O}(N/\lambda) \text{ iter.})$$

$$2. z_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \text{ for } i = 1, \dots, \lambda$$

$$3. x_i = \mathbf{m}^{(t)} + \sigma^{(t)} \sqrt{\mathbf{C}^{(t)}} z_i$$

$$4. (x_{i:\lambda})_{i=1,\dots,\lambda} = \text{SORTW.R.T.}^f((x_i)_{i=1,\dots,\lambda})$$

$$5. \mathbf{p}_c^{(t+1)} = (1 - c_c) \mathbf{p}_c^{(t)} + \sqrt{c_c(2 - c_c)/(\sum w_i^2)} \sum_{i=1}^{\lambda} w_i (x_{i:\lambda} - \mathbf{m}^{(t)}) / \sigma^{(t)}$$

$$6. \mathbf{p}_{\sigma}^{(t+1)} = (1 - c_{\sigma}) \mathbf{p}_{\sigma}^{(t)} + \sqrt{c_{\sigma}(2 - c_{\sigma})/(\sum w_i^2)} \sum_{i=1}^{\lambda} w_i z_{i:\lambda}$$

$$7. \sigma^{(t+1)} = \sigma^{(t)} \exp \left(\frac{c_{\sigma}}{d_{\sigma}} \left(\frac{\|\mathbf{p}_{\sigma}^{(t+1)}\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)$$

$$8. \mathbf{m}^{(t+1)} = \mathbf{m}^{(t)} + \sum_{i=1}^{\lambda} w_i (x_{i:\lambda} - \mathbf{m}^{(t)})$$

$$9. \mathbf{C}^{(t+1)} = \mathbf{C}^{(t)} + c_{\mu} \sum_{i=1}^{\lambda} w_i \left(\frac{x_{i:\lambda} - \mathbf{m}^{(t)}}{\sigma^{(t)}} \frac{(x_{i:\lambda} - \mathbf{m}^{(t)})^T}{\sigma^{(t)}} - \mathbf{C}^{(t)} \right) \\ + c_1 (\mathbf{p}_c^{(t+1)} \mathbf{p}_c^{(t+1)T} - \mathbf{C}^{(t)})$$

$\mathcal{O}(\lambda N^2)$

$\mathcal{O}(\lambda N)$

$\mathcal{O}(\lambda N^2)$

113

Sep-CMA-ES as a Natural Gradient

$\mathcal{O}(\lambda N)$ floating point (FP) multiplication + $\mathcal{O}(N + \lambda N)$ FP memory. ($\lambda \in o(N)$)

$$1. \sqrt{\mathbf{C}^{(t)}} = \text{MATRIXSQRT}(\mathbf{C}^{(t)})$$

$$2. z_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \text{ for } i = 1, \dots, \lambda$$

$$3. x_i = \mathbf{m}^{(t)} + \sigma^{(t)} \sqrt{\mathbf{C}^{(t)}} z_i$$

$$4. (x_{i:\lambda})_{i=1,\dots,\lambda} = \text{SORTW.R.T.}^f((x_i)_{i=1,\dots,\lambda})$$

$$5. \mathbf{p}_c^{(t+1)} = (1 - c_c) \mathbf{p}_c^{(t)} + \sqrt{c_c(2 - c_c)/(\sum w_i^2)} \sum_{i=1}^{\lambda} w_i (x_{i:\lambda} - \mathbf{m}^{(t)}) / \sigma^{(t)}$$

$$6. \mathbf{p}_{\sigma}^{(t+1)} = (1 - c_{\sigma}) \mathbf{p}_{\sigma}^{(t)} + \sqrt{c_{\sigma}(2 - c_{\sigma})/(\sum w_i^2)} \sum_{i=1}^{\lambda} w_i z_{i:\lambda}$$

$$7. \sigma^{(t+1)} = \sigma^{(t)} \exp \left(\frac{c_{\sigma}}{d_{\sigma}} \left(\frac{\|\mathbf{p}_{\sigma}^{(t+1)}\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)$$

$$8. \mathbf{m}^{(t+1)} = \mathbf{m}^{(t)} + \sum_{i=1}^{\lambda} w_i (x_{i:\lambda} - \mathbf{m}^{(t)})$$

$$9. \mathbf{C}^{(t+1)} = \mathbf{C}^{(t)} + c_{\mu} \sum_{i=1}^{\lambda} w_i \text{diag} \left(\frac{x_{i:\lambda} - \mathbf{m}^{(t)}}{\sigma^{(t)}} \frac{(x_{i:\lambda} - \mathbf{m}^{(t)})^T}{\sigma^{(t)}} - \mathbf{C}^{(t)} \right) \\ + c_1 \text{diag}(\mathbf{p}_c^{(t+1)} \mathbf{p}_c^{(t+1)T} - \mathbf{C}^{(t)})$$

$\mathcal{O}(\lambda N)$

$\mathcal{O}(\lambda N)$

$\mathcal{O}(\lambda N)$

Plug diagonal \mathbf{C} and compute the natural gradient

114

The Continuous Search Problem

Difficulties of a non-linear optimization problem are

- dimensionality and non-separability
demands to exploit problem structure, e.g. neighborhood
cave: design of benchmark functions
- ill-conditioning
demands to acquire a second order model
- ruggedness
demands a non-local (stochastic? population based?) approach

- 1 Problem Statement
- 2 Evolution Strategy (ES)
- 3 Step-Size Adaptation
 - ▶ Why Step-Size Control
 - ▶ Path Length Control (CSA)
 - ▶ Limitations of CSA and its Alternatives
- 4 Covariance Matrix Adaptation (CMA)
 - ▶ Rank-One Update and Cumulation
 - ▶ Rank- μ Update
 - ▶ Active Covariance Update
- 5 Design Principle
 - ▶ Theoretical Foundations
 - ▶ Variants for Large Scale Problems
- 6 Summary and Final Remarks

115

561

116

Main Characteristics of (CMA) Evolution Strategies

- ➊ Multivariate normal distribution to generate new search points follows the maximum entropy principle
- ➋ Rank-based selection implies invariance, same performance on $g(f(x))$ for any increasing g more invariance properties are featured
- ➌ Step-size control facilitates fast (log-linear) convergence and possibly linear scaling with the dimension in CMA-ES based on an evolution path (a non-local trajectory)
- ➍ Covariance matrix adaptation (CMA) increases the likelihood of previously successful steps and can improve performance by orders of magnitude

the update follows the natural gradient
 $\mathbf{C} \propto \mathbf{H}^{-1} \iff$ adapts a variable metric
 \iff new (rotated) problem representation
 $\implies f : \mathbf{x} \mapsto g(\mathbf{x}^T \mathbf{H} \mathbf{x})$ reduces to $\mathbf{x} \mapsto \mathbf{x}^T \mathbf{x}$

117



Thank You

Source code for CMA-ES in C, Java, Matlab, Octave, Python, Scilab is available at http://www.lri.fr/~hansen/cmaes_inmatlab.html

119



Limitations

of CMA Evolution Strategies

- ➊ internal CPU-time: $10^{-8}n^2$ seconds per function evaluation on a 2GHz PC, tweaks are available
 1 000 000 f -evaluations in 100-D take 100 seconds internal CPU-time
- ➋ better methods are presumably available in case of
 - ▶ partly separable problems
 - ▶ specific problems, for example with cheap gradients
 - ▶ small dimension ($n \ll 10$)
 for example Nelder-Mead
 - ▶ small running times (number of f -evaluations $< 100n$)
 model-based methods

118

